

Richard Blum and Christine Bresnahan

Linux[®] Command Line and Shell Scripting

Second Edition

Use command lines and
bypass the GUI

Automate common tasks

Create professional,
real-world scripts



Bible

The book you need to succeed!

Table of Contents

[*Title Page*](#)

[*Copyright*](#)

[*Dedication*](#)

[*Credits*](#)

[*About the Authors*](#)

[*About the Technical Editor*](#)

[*Acknowledgments*](#)

[*Introduction*](#)

[*Who Should Read This Book*](#)

[*How This Book Is Organized*](#)

[*Minimum Requirements*](#)

[*Where to Go from Here*](#)

Part I: The Linux Command Line

Chapter 1: Starting with Linux Shells

What Is Linux?

Linux Distributions

Summary

Chapter 2: Getting to the Shell

Terminal Emulation

The terminfo Database

The Linux Console

The xterm Terminal

The Konsole Terminal

The GNOME Terminal

Summary

Chapter 3: Basic bash Shell Commands

Starting the Shell

[*The Shell Prompt*](#)

[*The bash Manual*](#)

[*Filesystem Navigation*](#)

[*File and Directory Listing*](#)

[*File Handling*](#)

[*Directory Handling*](#)

[*Viewing File Contents*](#)

[*Summary*](#)

[*Chapter 4: More bash Shell Commands*](#)

[*Monitoring Programs*](#)

[*Monitoring Disk Space*](#)

[*Working with Data Files*](#)

[*Summary*](#)

[*Chapter 5: Using Linux Environment Variables*](#)

[*What Are Environment Variables?*](#)

[*Setting Environment Variables*](#)

[*Removing Environment Variables*](#)

[**Default Shell Environment**](#)

[**Variables**](#)

[**Setting the PATH Environment**](#)

[**Variable**](#)

[**Locating System Environment**](#)

[**Variables**](#)

[**Variable Arrays**](#)

[**Using Command Aliases**](#)

[**Summary**](#)

[**Chapter 6: Understanding Linux File Permissions**](#)

[**Linux Security**](#)

[**Using Linux Groups**](#)

[**Decoding File Permissions**](#)

[**Changing Security Settings**](#)

[**Sharing Files**](#)

[**Summary**](#)

[**Chapter 7: Managing Filesystems**](#)

[Exploring Linux Filesystems](#)
[Working with Filesystems](#)
[Logical Volume Managers](#)
[Summary](#)

[Chapter 8: Installing Software](#)

[Package Management Primer](#)
[The Debian-Based Systems](#)
[The Red Hat-Based Systems](#)
[Installing from Source Code](#)
[Summary](#)

[Chapter 9: Working with Editors](#)

[The vim Editor](#)
[The emacs Editor](#)
[The KDE Family of Editors](#)
[The GNOME Editor](#)
[Summary](#)

[Part II: Shell Scripting Basics](#)

Chapter 10: Basic Script Building

Using Multiple Commands

Creating a Script File

Displaying Messages

Using Variables

Redirecting Input and Output

Pipes

Performing Math

Exiting the Script

Summary

Chapter 11: Using Structured Commands

Working with the if-then Statement

The if-then-else Statement

Nesting ifs

The test Command

Compound Condition Testing

Advanced if-then Features

[The case Command](#) [Summary](#)

[Chapter 12: More Structured Commands](#)

[The for Command](#)

[The C-Style for Command](#)

[The while Command](#)

[The until Command](#)

[Nesting Loops](#)

[Looping on File Data](#)

[Controlling the Loop](#)

[Processing the Output of a Loop](#)

[Summary](#)

[Chapter 13: Handling User Input](#)

[Command Line Parameters](#)

[Special Parameter Variables](#)

[Being Shifty](#)

[Working with Options](#)

[Standardizing Options](#)
[Getting User Input](#)
[Summary](#)

[Chapter 14: Presenting Data](#)

[Understanding Input and Output](#)
[Redirecting Output in Scripts](#)
[Redirecting Input in Scripts](#)
[Creating Your Own Redirection](#)
[Listing Open File Descriptors](#)
[Suppressing Command Output](#)
[Using Temporary Files](#)
[Logging Messages](#)
[Summary](#)

[Chapter 15: Script Control](#)

[Handling Signals](#)
[Running Scripts in Background Mode](#)
[Running Scripts Without a Console](#)

[**Job Control**](#)

[**Being Nice**](#)

[**Running Like Clockwork**](#)

[**Summary**](#)

[**Part III: Advanced Shell Scripting**](#)

[**Chapter 16: Creating Functions**](#)

[**Basic Script Functions**](#)

[**Returning a Value**](#)

[**Using Variables in Functions**](#)

[**Array Variables and Functions**](#)

[**Function Recursion**](#)

[**Creating a Library**](#)

[**Using Functions on the Command Line**](#)

[**Summary**](#)

[**Chapter 17: Writing Scripts for**](#)

Graphical Desktops

Creating Text Menus

Doing Windows

Getting Graphic

Summary

Chapter 18: Introducing sed and gawk

Text Manipulation

The sed Editor Basics

Summary

Chapter 19: Regular Expressions

What Are Regular Expressions?

Defining BRE Patterns

Extended Regular Expressions

Regular Expressions in Action

Summary

Chapter 20: Advanced sed

Multiline Commands

[The Hold Space](#)

[Negating a Command](#)

[Changing the Flow](#)

[Pattern Replacement](#)

[Using sed in Scripts](#)

[Creating sed Utilities](#)

[Summary](#)

[Chapter 21: Advanced gawk](#)

[Using Variables](#)

[Working with Arrays](#)

[Using Patterns](#)

[Structured Commands](#)

[Formatted Printing](#)

[Built-in Functions](#)

[User-Defined Functions](#)

[Summary](#)

[Chapter 22: Working with Alternative Shells](#)

[**What Is the dash Shell?**](#)

[**The dash Shell Features**](#)

[**Scripting in dash**](#)

[**The zsh Shell**](#)

[**Parts of the zsh Shell**](#)

[**Scripting with zsh**](#)

[**Summary**](#)

[**Part IV: Advanced Shell Scripting Topics**](#)

[**Chapter 23: Using a Database**](#)

[**The MySQL Database**](#)

[**The PostgreSQL Database**](#)

[**Working with Tables**](#)

[**Using the Database in Your Scripts**](#)

[**Summary**](#)

[**Chapter 24: Using the Web**](#)

[**The Lynx Program**](#)

[The cURL Program](#)
[Networking with zsh](#)
[Summary](#)

[Chapter 25: Using E-mail](#)

[The Basics of Linux E-Mail](#)
[Setting Up Your Server](#)
[Sending a Message with Mailx](#)
[The Mutt Program](#)
[Summary](#)

[Chapter 26: Writing Script Utilities](#)

[Monitoring Disk Space](#)
[Performing Backups](#)
[Managing User Accounts](#)
[Summary](#)

[Chapter 27: Advanced Shell Scripts](#)

[Monitoring System Statistics](#)
[Problem-Tracking Database](#)

[Summary](#)

[Appendix A: Quick Guide to bash Commands](#)

[Built-in Commands](#)

[bash Commands](#)

[Environment Variables](#)

[Appendix B: Quick Guide to sed and gawk](#)

[The sed Editor](#)

[The gawk Program](#)

[Index](#)

Linux[®] Command Line and Shell Scripting Bible

Second Edition

**Richard Blum
Christine Bresnahan**



Wiley Publishing, Inc.

Linux® Command Line and Shell Scripting Bible, Second Edition

Published by

Wiley Publishing, Inc.

10475 Crosspoint Boulevard

Indianapolis, IN 46256

www.wiley.com

Copyright © 2011 by Wiley Publishing, Inc., Indianapolis,
Indiana

Published simultaneously in Canada

ISBN: 978-1-118-00442-5

ISBN: 978-1-118-08777-0 (ebk)

ISBN: 978-1-118-08778-7 (ebk)

ISBN: 978-1-118-08779-4 (ebk)

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 646-

8600. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permissions>.

Limit of Liability/Disclaimer of Warranty: The publisher and the author make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation warranties of fitness for a particular purpose. No warranty may be created or extended by sales or promotional materials. The advice and strategies contained herein may not be suitable for every situation. This work is sold with the understanding that the publisher is not engaged in rendering legal, accounting, or other professional services. If professional assistance is required, the services of a competent professional person should be sought. Neither the publisher nor the author shall be liable for damages arising herefrom. The fact that an organization or Web site is referred to in this work as a citation and/or a potential source of further information does not mean that the author or the publisher endorses the information the organization or website may provide or recommendations it may make. Further, readers should be aware that Internet websites listed in this work may have changed or disappeared between when this work was written and when it is read.

For general information on our other products and services

please contact our Customer Care Department within the United States at (877) 762-2974, outside the United States at (317) 572-3993 or fax (317) 572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic books.

Library of Congress Control Number: 2011921770

Trademarks: Wiley and the Wiley logo are trademarks or registered trademarks of John Wiley & Sons, Inc. and/or its affiliates, in the United States and other countries, and may not be used without written permission. Linux is a registered trademark of Linus Torvalds. All other trademarks are the property of their respective owners. Wiley Publishing, Inc. is not associated with any product or vendor mentioned in this book.

To the Lord God Almighty, “in whom are hidden all the treasures of wisdom and knowledge.”

—Colossians 2:3

Credits

Acquisitions Editor

Mary James

Project Editor

Brian Herrmann

Technical Editor

Jack Cox

Production Editor

Daniel Scribner

Copy Editor

Nancy Rapoport

Editorial Director

Robyn B. Siesky

Editorial Manager

Mary Beth Wakefield

Freelancer Editorial Manager

Rosemarie Graham

Marketing Manager

Ashley Zurcher

Production Manager

Tim Tate

Vice President and Executive Group Publisher

Richard Swadley

Vice President and Executive Publisher

Barry Pruett

Associate Publisher

Jim Minatel

Project Coordinator, Cover

Katherine Crocker

Proofreader

Publication Services, Inc.

Indexer

Robert Swanson

Cover Designer

Ryan Sneed

Cover Image

Joyce Haughey

About the Authors

Richard Blum has worked in the IT industry for more than 20 years as both a systems and network administrator, and has published numerous Linux and open source books. He has administered UNIX, Linux, Novell, and Microsoft servers, as well as helped design and maintain a 3,500 user network utilizing Cisco switches and routers. He has used Linux servers and shell scripts to perform automated network monitoring, and has written shell scripts in most of the common Linux shell environments. Rich is an online instructor for an Introduction to Linux course that is used by colleges and universities across the U.S. When he's not being a computer nerd, Rich plays electric bass in a couple of different church worship bands, and enjoys spending time with his wife Barbara, and two daughters Katie Jane and Jessica.

Christine Bresnahan starting working with computers more than 25 years ago in the IT industry as a system administrator. Christine is currently an Adjunct Professor at Ivy Tech Community College in Indianapolis, Indiana, teaching Linux system administration, Linux security, and Windows security classes.

About the Technical Editor

Jack Cox is a Senior Manager with CapTech Ventures in Richmond, VA. He has more than 25 years of experience in IT, spanning a broad range of disciplines including mobile computing, transaction processing, RFID, Java development, and cryptography. Jack enjoys life in Richmond with his lovely wife and rambunctious children. Outside of technology, his interests include church, his children, and extended family.

Acknowledgments

First, all glory and praise go to God, who through His Son, Jesus Christ, makes all things possible, and gives us the gift of eternal life.

Many thanks go to the fantastic team of people at John Wiley & Sons for their outstanding work on this project. Thanks to Mary James, the acquisitions editor, for offering us the opportunity to work on this book. Also thanks to Brian Herrmann, the project editor, for keeping things on track and making this book more presentable. Thanks, Brian, for all your hard work and diligence. The technical editor, Jack Cox, did a wonderful job of double checking all the work in the book, plus making suggestions to improve the content. Thanks to Nancy Rapoport, the copy editor, for her endless patience and diligence to make our work readable. We would also like to thank Carole McClendon at Waterside Productions, Inc., for arranging this opportunity for us, and for helping us out in our writing careers.

Christine would like to thank her husband, Timothy, for his encouragement, patience, and willingness to listen, even when he has no idea what she is talking about.

Introduction

Welcome to the second edition of *Linux Command Line and Shell Scripting Bible*. Like all books in the *Bible* series, you can expect to find both hands-on tutorials and real-world information, as well as reference and background information that provides a context for what you are learning. This book is a fairly comprehensive resource on the Linux command line and shell commands. By the time you have completed *Linux Command Line and Shell Scripting Bible* you will be well prepared to write your own shell scripts that can automate practically any task on your Linux system.

Who Should Read This Book

If you're a system administrator in a Linux environment, you'll benefit greatly by knowing how to write shell scripts. The book doesn't walk you through the process of setting up a Linux system, but once you have it running you'll want to start automating some of the routine administrative tasks. That's where shell scripting comes in, and that's where this book will help you out. This book will demonstrate how to automate any administrative task using shell scripts, from monitoring system statistics and data

files to generating reports for your boss.

If you're a home Linux enthusiast, you'll also benefit from *Linux Command Line and Shell Scripting Bible*. Nowadays, it's easy to get lost in the graphical world of pre-built widgets. Most desktop Linux distributions try their best to hide the Linux system from the typical user. However, there are times when you have to know what's going on under the hood. This book shows you how to access the Linux command line prompt and what to do once you get there. Often performing simple tasks, such as file management, can be done more quickly from the command line than from a fancy graphical interface. There's a wealth of commands you can use from the command line, and this book shows you how to use them.

How This Book Is Organized

This book leads you through the basics of the Linux command line and into more complicated topics, such as creating your own shell scripts. The book is divided into five parts, each one building on the previous parts.

Part I assumes that you either have a Linux system running, or are looking into getting a Linux system. Chapter 1, "Starting with Linux Shells," describes the parts of a total Linux system and shows how the shell fits in. After describing the basics of the Linux system, this part continues with the following:

- Using a terminal emulation package to access the

shell (Chapter 2)

- Introducing the basic shell commands (Chapter 3)
- Using more advanced shell commands to peek at system information (Chapter 4)
- Working with shell variables to manipulate data (Chapter 5)
- Understanding the Linux filesystem and security (Chapter 6)
- Working with Linux filesystems from the command line (Chapter 7)
- Installing and updating software from the command line (Chapter 8)
- Using the Linux editors to start writing shell scripts (Chapter 9)

In Part II, you begin writing shell scripts. As you go through the chapters you'll do the following:

- Learn how to create and run shell scripts (Chapter 10)
- Alter the program flow in a shell script (Chapter 11)
- Iterate through code sections (Chapter 12)
- Handle data from the user in your scripts (Chapter 13)
- See different methods for storing and displaying data from your Script (Chapter 14)
- Control how and when your shell scripts run on the system (Chapter 15)

Part III dives into more advanced areas of shell script

programming, including:

- Creating your own functions to use in all your scripts (Chapter 16)
- Utilizing the Linux graphical desktop for interacting with your script users (Chapter 17)
- Using advanced Linux commands to filter and parse data files (Chapter 18)
- Using regular expressions to define data (Chapter 19)
- Learning advanced methods of manipulating data in your scripts (Chapter 20)
- Generating reports from raw data (Chapter 21)
- Modifying your shell scripts to run in other Linux shells (Chapter 22)

The last section of the book, Part IV, demonstrates how to use shell scripts in real-world environments. In this part, you will:

- See how to use popular open source databases in your shell scripts (Chapter 23)
- Learn how to extract data from Websites, and send data between systems (Chapter 24)
- Use e-mail to send notifications and reports to external users (Chapter 25)
- Write shell scripts to automate your daily system administration functions (Chapter 26)
- Utilize all of the features you've learned from the book to create professional-quality shell scripts (Chapter

Conventions and Features

There are many different organizational and typographical features throughout this book designed to help you get the most of the information.

Notes and Warnings

Whenever the authors want to bring something important to your attention the information will appear in a Warning.

Warning

This information is important and is set off in a separate paragraph with a special icon. Warnings provide information about things to watch out for, whether simply inconvenient or potentially hazardous to your data or systems.

For additional items of interest that relate to the chapter text, the authors will use Notes.

Note

Notes provide additional, ancillary information that is helpful, but somewhat outside of the current presentation of information.

Minimum Requirements

Linux Command Line and Shell Scripting Bible doesn't focus on any specific Linux distribution, so you'll be able to follow along in the book using any Linux system you have available. The bulk of the book references the bash shell, which is the default shell for most Linux systems.

Where to Go from Here

Once you've completed *Linux Command Line and Shell Scripting Bible*, you'll be well on your way to incorporating Linux commands in your daily Linux work. In the ever-changing world of Linux, it's always a good idea to stay in touch with new developments. Often Linux distributions will change, adding new features and removing older ones. To keep your knowledge of Linux fresh, always stay well-informed. Find a good Linux forum site and monitor what's happening in the Linux world. There are many popular Linux news sites, such as Slashdot and Distrowatch, that provide up-to-the-minute information about new advances in Linux.

Part I

The Linux Command Line

In This Part

Chapter 1: Starting with Linux Shells

Chapter 2: Getting to the Shell

Chapter 3: Basic bash Shell Commands

Chapter 4: More bash Shell Commands

Chapter 5: Using Linux Environment Variables

Chapter 6: Understanding Linux File Permissions

Chapter 7: Managing Filesystems

Chapter 8: Installing Software

Chapter 9: Working with Editors

Chapter 1

Starting with Linux Shells

In this Chapter

- What is Linux?
- Parts of the Linux kernel
- Exploring the Linux desktop
- Visiting Linux distributions

Before you can dive into working with the Linux command line and shells, it's a good idea to first understand what Linux is, where it came from, and how it works. This chapter walks you through what Linux is, and explains where the shell and command line fit in the overall Linux picture.

What Is Linux?

If you've never worked with Linux before, you may be confused as to why there are so many different versions of it available. I'm sure that you have heard various terms such as distribution, LiveCD, and GNU when looking at Linux packages and been confused. Wading through the world of

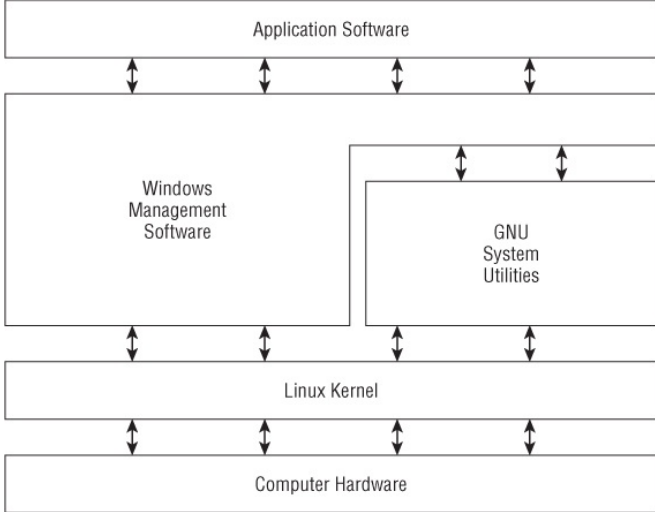
Linux for the first time can be a tricky experience. This chapter takes some of the mystery out of the Linux system before you start working on commands and scripts.

For starters, there are four main parts that make up a Linux system:

- The Linux kernel
- The GNU utilities
- A graphical desktop environment
- Application software

Each of these four parts has a specific job in the Linux system. Each of the parts by itself isn't very useful. [Figure 1.1](#) shows a basic diagram of how the parts fit together to create the overall Linux system.

[Figure 1.1](#) The Linux system



This section describes these four main parts in detail, and gives you an overview of how they work together to create a complete Linux system.

Looking into the Linux Kernel

The core of the Linux system is the *kernel*. The kernel controls all of the hardware and software on the computer system, allocating hardware when necessary, and executing software when required.

If you've been following the Linux world at all, no doubt

you've heard the name Linus Torvalds. Linus is the person responsible for creating the first Linux kernel software while he was a student at the University of Helsinki. He intended it to be a copy of the Unix system, at the time a popular operating system used at many universities.

After developing the Linux kernel, Linus released it to the Internet community and solicited suggestions for improving it. This simple process started a revolution in the world of computer operating systems. Soon Linus was receiving suggestions from students as well as professional programmers from around the world.

Allowing anyone to change programming code in the kernel would result in complete chaos. To simplify things, Linus acted as a central point for all improvement suggestions. It was ultimately Linus's decision whether or not to incorporate suggested code in the kernel. This same concept is still in place with the Linux kernel code, except that instead of just Linus controlling the kernel code, a team of developers has taken on the task.

The kernel is primarily responsible for four main functions:

- System memory management
- Software program management
- Hardware management
- Filesystem management

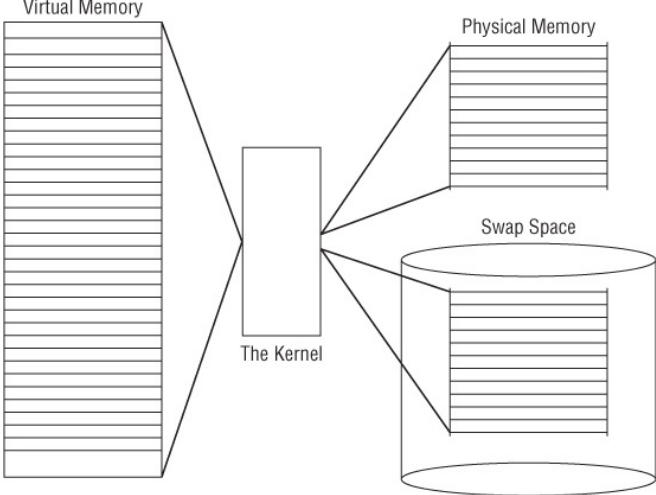
The following sections explore each of these functions in more detail.

System Memory Management

One of the primary functions of the operating system kernel is memory management. Not only does the kernel manage the physical memory available on the server, but it can also create and manage virtual memory, or memory that does not actually exist.

It does this by using space on the hard disk, called the *swap space*. The kernel swaps the contents of virtual memory locations back and forth from the swap space to the actual physical memory. This allows the system to think there is more memory available than what physically exists (shown in [Figure 1.2](#)).

[Figure 1.2](#) The Linux system memory map



The memory locations are grouped into blocks called *pages*. The kernel locates each page of memory either in the physical memory or the swap space. The kernel then maintains a table of the memory pages that indicates which pages are in physical memory and which pages are swapped out to disk.

The kernel keeps track of which memory pages are in use and automatically copies memory pages that have not been accessed for a period of time to the swap space area (called *swapping out*), even if there's other memory available. When a program wants to access a memory

page that has been swapped out, the kernel must make room for it in physical memory by swapping out a different memory page, and swap in the required page from the swap space. Obviously, this process takes time, and can slow down a running process. The process of swapping out memory pages for running applications continues for as long as the Linux system is running.

You can see the current status of the virtual memory on your Linux system by viewing the special `/proc/meminfo` file. Here's an example of a sample `/proc/meminfo` entry:

```
rich@rich-desktop:~$ cat /proc/meminfo
MemTotal:        1026084 kB
MemFree:         666356 kB
Buffers:         49900 kB
Cached:          152272 kB
SwapCached:       0 kB
Active:          171468 kB
Inactive:        154196 kB
Active(anon):    131056 kB
Inactive(anon):   32 kB
Active(file):    40412 kB
Inactive(file):  154164 kB
Unevictable:     12 kB
Mlocked:         12 kB
HighTotal:       139208 kB
HighFree:        252 kB
LowTotal:        886876 kB
LowFree:         666104 kB
SwapTotal:       2781176 kB
SwapFree:        2781176 kB
Dirty:           588 kB
Writeback:       0 kB
AnonPages:       123500 kB
Mapped:          52232 kB
Shmem:           7600 kB
Slab:            17676 kB
SReclaimable:    9788 kB
SUnreclaim:     7888 kB
KernelStack:    2656 kB
PageTables:      5072 kB
NFS_Unstable:    0 kB
Bounce:          0 kB
```

```
WritebackTmp:      0 kB
CommitLimit:      3294216 kB
Committed_AS:     1234480 kB
VmallocTotal:     122880 kB
VmallocUsed:       7520 kB
VmallocChunk:     110672 kB
HardwareCorrupted: 0 kB
HugePages_Total:   0
HugePages_Free:    0
HugePages_Rsvd:    0
HugePages_Surp:    0
Hugepagesize:      4096 kB
DirectMap4k:       12280 kB
DirectMap4M:       897024 kB
rich@rich-desktop:~$
```

The `MemTotal:` line shows that this Linux server has 1GB of physical memory. It also shows that about 660MB is not currently being used (`MemFree`). The output also shows that there is about 2.5GB of swap space memory available on this system (`SwapTotal`).

By default, each process running on the Linux system has its own private memory pages. One process cannot access memory pages being used by another process. The kernel maintains its own memory areas. For security purposes, no processes can access memory used by the kernel processes.

To facilitate data sharing, you can create shared memory pages. Multiple processes can read and write to and from a common shared memory area. The kernel maintains and administers the shared memory areas and allows individual processes access to the shared area.

The special `ipcs` command allows you to view the current shared memory pages on the system. Here's the output from a sample `ipcs` command:

```
# ipcs -m
```

----- Shared Memory Segments -----

key	shmid	owner	perms	bytes	nattch	status
0x00000000	0	rich	600	52228	6	dest
0x395ec51c	1	oracle	640	5787648	6	

#

Each shared memory segment has an owner that created the segment. Each segment also has a standard Linux permissions setting that sets the availability of the segment for other users. The key value is used to allow other users to gain access to the shared memory segment.

Software Program Management

The Linux operating system calls a running program a *process*. A process can run in the foreground, displaying output on a display, or it can run in background, behind the scenes. The kernel controls how the Linux system manages all the processes running on the system.

The kernel creates the first process, called the *init process*, to start all other processes on the system. When the kernel starts, it loads the init process into virtual memory. As the kernel starts each additional process, it gives it a unique area in virtual memory to store the data and code that the process uses.

Some Linux implementations contain a table of processes to start automatically on bootup. On Linux systems, this table is usually located in the special file

`/etc/inittabs`.

Other systems (such as the popular Ubuntu Linux

distribution) utilize the `/etc/init.d` folder, which contains scripts for starting and stopping individual applications at boot time. The scripts are started via entries under the `/etc/rcX.d` folders, where *X* is a *run level*.

The Linux operating system uses an init system that utilizes run levels. A run level can be used to direct the init process to run only certain types of processes, as defined in the `/etc/inittabs` file or the `/etc/rcX.d` folders. There are five init run levels in the Linux operating system.

At run level 1, only the basic system processes are started, along with one console terminal process. This is called *single user* mode. Single user mode is most often used for emergency filesystem maintenance when something is broken. Obviously, in this mode only one person (usually the administrator) can log in to the system to manipulate data.

The standard init run level is 3. At this run level, most application software such as network support software is started. Another popular run level in Linux is run level 5. This is the run level where the system starts the graphical X Window software, and allows you to log in using a graphical desktop window.

The Linux system can control the overall system functionality by controlling the init run level. By changing the run level from 3 to 5, the system can change from a console-based system to an advanced, graphical X Window system.

In Chapter 4, you'll see how to use the `ps` command to

view the processes currently running on the Linux system. Here's an example of what you'll see using the `ps` command:

```
$ ps ax
  PID TTY          STAT       TIME COMMAND
    1 ?            S           0:03  init
    2 ?            SW          0:00  [kflushd]
    3 ?            SW          0:00  [kupdate]
    4 ?            SW          0:00  [kpiod]
    5 ?            SW          0:00  [kswapd]
  243 ?            SW          0:00  [portmap]
  295 ?            S           0:00  syslogd
  305 ?            S           0:00  klogd
  320 ?            S           0:00  /usr/sbin/atd
  335 ?            S           0:00  crond
  350 ?            S           0:00  inetd
  365 ?            SW          0:00  [lpd]
  403 ttyS0        S           0:00  gpm -t ms
  418 ?            S           0:00  httpd
  423 ?            S           0:00  httpd
  424 ?            SW          0:00  [httpd]
  425 ?            SW          0:00  [httpd]
  426 ?            SW          0:00  [httpd]
  427 ?            SW          0:00  [httpd]
  428 ?            SW          0:00  [httpd]
  429 ?            SW          0:00  [httpd]
  430 ?            SW          0:00  [httpd]
  436 ?            SW          0:00  [httpd]
  437 ?            SW          0:00  [httpd]
  438 ?            SW          0:00  [httpd]
  470 ?            S           0:02  xfs -port -l
  485 ?            SW          0:00  [smbd]
  495 ?            S           0:00  nmbd -D
  533 ?            SW          0:00  [postmaster]
  538 tty1        SW          0:00  [mingetty]
  539 tty2        SW          0:00  [mingetty]
  540 tty3        SW          0:00  [mingetty]
  541 tty4        SW          0:00  [mingetty]
  542 tty5        SW          0:00  [mingetty]
  543 tty6        SW          0:00  [mingetty]
  544 ?            SW          0:00  [prefdm]
  549 ?            SW          0:00  [prefdm]
  559 ?            S           0:02  [kwm]
  585 ?            S           0:06  kikbd
  594 ?            S           0:00  kwmsound
  595 ?            S           0:03  kpanel
  596 ?            S           0:02  kfm
  597 ?            S           0:00  krootwm
  598 ?            S           0:01  kbgndwm
```



```

611 ?      S      0:00 kcmllaptop -daemon
666 ?      S      0:00 /usr/libexec/postfix/master
668 ?      S      0:00 qmgr -l -t fifo -u
787 ?      S      0:00 pickup -l -t fifo
790 ?      S      0:00 telnetd: 192.168.1.2 [vt100]
791 pts/0   S      0:00 login -- rich
792 pts/0   S      0:00 -bash
805 pts/0   R      0:00 ps ax

```

\$

The first column in the output shows the *process ID* (or PID) of the process. Notice that the first process is our friend the init process, and assigned PID 1 by the Linux system. All other processes that start after the init process are assigned PIDs in numerical order. No two processes can have the same PID (although old PID numbers can be reused by the system after the original process terminates).

The third column shows the current status of the process (*s* for sleeping, *sw* for sleeping and waiting, and *R* for running). The process name is shown in the last column. Processes that are in brackets are processes that have been swapped out of memory to the disk swap space due to inactivity. You can see that some of the processes have been swapped out, but most of the running processes have not.

Hardware Management

Still another responsibility for the kernel is hardware management. Any device that the Linux system must communicate with needs driver code inserted inside the kernel code. The driver code allows the kernel to pass data back and forth to the device, acting as a middle man

between applications and the hardware. There are two methods used for inserting device driver code in the Linux kernel:

- Drivers compiled in the kernel
- Driver modules added to the kernel

Previously, the only way to insert device driver code was to recompile the kernel. Each time you added a new device to the system, you had to recompile the kernel code. This process became even more inefficient as Linux kernels supported more hardware. Fortunately, Linux developers devised a better method to insert driver code into the running kernel.

Programmers developed the concept of kernel modules to allow you to insert driver code into a running kernel without having to recompile the kernel. Also, a kernel module could be removed from the kernel when the device was finished being used. This greatly simplified and expanded using hardware with Linux.

The Linux system identifies hardware devices as special files, called *device files*. There are three different classifications of device files:

- Character
- Block
- Network

Character device files are for devices that can only handle data one character at a time. Most types of modems and terminals are created as character files. Block files are for devices that can handle data in large

blocks at a time, such as disk drives.

The network file types are used for devices that use packets to send and receive data. This includes network cards and a special loopback device that allows the Linux system to communicate with itself using common network programming protocols.

Linux creates special files, called nodes, for each device on the system. All communication with the device is performed through the device node. Each node has a unique number pair that identifies it to the Linux kernel. The number pair includes a major and a minor device number. Similar devices are grouped into the same major device number. The minor device number is used to identify a specific device within the major device group. The following is an example of a few device files on a Linux server:

```
rich@rich-desktop: ~$ cd /dev
rich@rich-desktop:/dev$ ls -al sda* ttyS*
brw-rw---- 1 root disk      8,  0 2010-09-18 17:25 sda
brw-rw---- 1 root disk      8,  1 2010-09-18 17:25 sda1
brw-rw---- 1 root disk      8,  2 2010-09-18 17:25 sda2
brw-rw---- 1 root disk      8,  5 2010-09-18 17:25 sda5
crw-rw---- 1 root dialout  4, 64 2010-09-18 17:25 ttyS0
crw-rw---- 1 root dialout  4, 65 2010-09-18 17:25 ttyS1
crw-rw---- 1 root dialout  4, 66 2010-09-18 17:25 ttyS2
crw-rw---- 1 root dialout  4, 67 2010-09-18 17:25 ttyS3
rich@rich-desktop:/dev$
```

Different Linux distributions handle devices using different device names. In this distribution, the `sda` device is the first ATA hard drive, and the `ttyS` devices are the standard IBM PC COM ports. The listing shows all of the `sda` devices that were created on the sample Linux system. Not all are actually used, but they are created in case the

administrator needs them. Similarly, the listing shows all of the `ttyS` devices created.

The fifth column is the major device node number. Notice that all of the `sda` devices have the same major device node, 8, while all of the `ttyS` devices use 4. The sixth column is the minor device node number. Each device within a major number has its own unique minor device node number.

The first column indicates the permissions for the device file. The first character of the permissions indicates the type of file. Notice that the ATA hard drive files are all marked as block (`b`) device, while the COM port device files are marked as character (`c`) devices.

Filesystem Management

Unlike some other operating systems, the Linux kernel can support different types of filesystems to read and write data to and from hard drives. Besides having over a dozen filesystems of its own, Linux can read and write to and from filesystems used by other operating systems, such as Microsoft Windows. The kernel must be compiled with support for all types of filesystems that the system will use. [Table 1.1](#) lists the standard filesystems that a Linux system can use to read and write data.

Table 1.1 Linux Filesystems

Filesystem	Description
ext	Linux Extended filesystem—the original Linux filesystem
ext2	Second extended filesystem, provided advanced features

ext2	over ext
ext3	Third extended filesystem, supports journaling
ext4	Fourth extended filesystem, supports advanced journaling
hpfs	OS/2 high-performance filesystem
jfs	IBMs journaling file system
iso9660	ISO 9660 filesystem (CD-ROMs)
minix	MINIX filesystem
msdos	Microsoft FAT16
ncp	Netware filesystem
nfs	Network File System
ntfs	Support for Microsoft NT filesystem
proc	Access to system information
ReiserFS	Advanced Linux file system for better performance and disk recovery
smb	Samba SMB filesystem for network access
sysv	Older Unix filesystem
ufs	BSD filesystem
umsdos	Unix-like filesystem that resides on top of msdos
vfat	Windows 95 filesystem (FAT32)
XFS	High-performance 64-bit journaling filesystem

Any hard drive that a Linux server accesses must be formatted using one of the filesystem types listed in [Table 1.1](#).

The Linux kernel interfaces with each filesystem using the Virtual File System (VFS). This provides a standard interface for the kernel to communicate with any type of filesystem. VFS caches information in memory as each

filesystem is mounted and used.

The GNU Utilities

Besides having a kernel to control hardware devices, a computer operating system needs utilities to perform standard functions, such as controlling files and programs. While Linus created the Linux system kernel, he had no system utilities to run on it. Fortunately for him, at the same time he was working, a group of people were working together on the Internet trying to develop a standard set of computer system utilities that mimicked the popular Unix operating system.

The GNU organization (GNU stands for GNU's Not Unix) developed a complete set of Unix utilities, but had no kernel system to run them on. These utilities were developed under a software philosophy called open source software (OSS).

The concept of OSS allows programmers to develop software and then release it to the world with no licensing fees attached. Anyone can use the software, modify it, or incorporate it into his or her own system without having to pay a license fee. Uniting Linus's Linux kernel with the GNU operating system utilities created a complete, functional, free operating system.

While the bundling of the Linux kernel and GNU utilities is often just called Linux, you will see some Linux purists on the Internet refer to it as the GNU/Linux system to give

credit to the GNU organization for its contributions to the cause.

The Core GNU Utilities

The GNU project was mainly designed for Unix system administrators to have a Unix-like environment available. This focus resulted in the project porting many common Unix system command line utilities. The core bundle of utilities supplied for Linux systems is called the *coreutils* package.

The GNU coreutils package consists of three parts:

- Utilities for handling files
- Utilities for manipulating text
- Utilities for managing processes

Each of these three main groups of utilities contains several utility programs that are invaluable to the Linux system administrator and programmer. This book covers each of the utilities contained in the GNU coreutils package in detail.

The Shell

The GNU/Linux shell is a special interactive utility. It provides a way for users to start programs, manage files on the filesystem, and manage processes running on the Linux system. The core of the shell is the command prompt. The command prompt is the interactive part of the shell. It allows you to enter text commands, and then it interprets the

commands and then executes them in the kernel.

The shell contains a set of internal commands that you use to control things such as copying files, moving files, renaming files, displaying the programs currently running on the system, and stopping programs running on the system. Besides the internal commands, the shell also allows you to enter the name of a program at the command prompt. The shell passes the program name off to the kernel to start it.

You can also group shell commands into files to execute as a program. Those files are called *shell scripts*. Any command that you can execute from the command line can be placed in a shell script and run as a group of commands. This provides great flexibility in creating utilities for commonly run commands, or processes that require several commands grouped together.

There are quite a few Linux shells available to use on a Linux system. Different shells have different characteristics, some being more useful for creating scripts and some being more useful for managing processes. The default shell used in all Linux distributions is the bash shell. The bash shell was developed by the GNU project as a replacement for the standard Unix shell, called the Bourne shell (after its creator). The bash shell name is a play on this wording, referred to as the “Bourne again shell.”

In addition to the bash shell, we will cover several other popular shells in this book. [Table 1.2](#) lists the different shells we will examine.

[Table 1.2](#) Linux Shells

Shell	Description
ash	A simple, lightweight shell that runs in low-memory environments but has full compatibility with the bash shell
korn	A programming shell compatible with the Bourne shell but supporting advanced programming features like associative arrays and floating-point arithmetic
tcsh	A shell that incorporates elements from the C programming language into shell scripts
zsh	An advanced shell that incorporates features from bash, tcsh, and korn, providing advanced programming features, shared history files, and themed prompts

Most Linux distributions include more than one shell, although usually they pick one of them to be the default. If your Linux distribution includes multiple shells, feel free to experiment with different shells and see which one fits your needs.

The Linux Desktop Environment

In the early days of Linux (the early 1990s) all that was available was a simple text interface to the Linux operating system. This text interface allowed administrators to start programs, control program operations, and move files around on the system.

With the popularity of Microsoft Windows, computer users expected more than the old text interface to work with. This spurred more development in the OSS community, and the Linux graphical desktops emerged.

Linux is famous for being able to do things in more than

one way, and no place is this more relevant than in graphical desktops. There are a plethora of graphical desktops you can choose from in Linux. The following sections describe a few of the more popular ones.

The X Windows System

There are two basic elements that control your video environment—the video card in your PC and your monitor. To display fancy graphics on your computer, the Linux software needs to know how to talk to both of them. The X Windows software is the core element in presenting graphics.

The X Windows software is a low-level program that works directly with the video card and monitor in the PC, and controls how Linux applications can present fancy windows and graphics on your computer.

Linux isn't the only operating system that uses X Windows; there are versions written for many different operating systems. In the Linux world, there are only two software packages that can implement it.

The XFree86 software package is the older of the two, and for a long time was the only X Windows package available for Linux. As its name implies, it's a free open source version of the X Windows software.

The newer of the two packages, X.org, has made great inroads in the Linux world and is now the more popular of the two. It, too, provides an open source software implementation of the X Windows system, but has support

for more of the newer video cards used today.

Both packages work the same way, controlling how Linux uses your video card to display content on your monitor. To do that, they have to be configured for your specific system. That is supposed to happen automatically when you install Linux.

When you first install a Linux distribution, it attempts to detect your video card and monitor, and then creates an X Windows configuration file that contains the required information. During installation you may notice a time when the installation program scans your monitor for supported video modes. Sometimes this causes your monitor to go blank for a few seconds. Because there are lots of different types of video cards and monitors out there, this process can take a little while to complete.

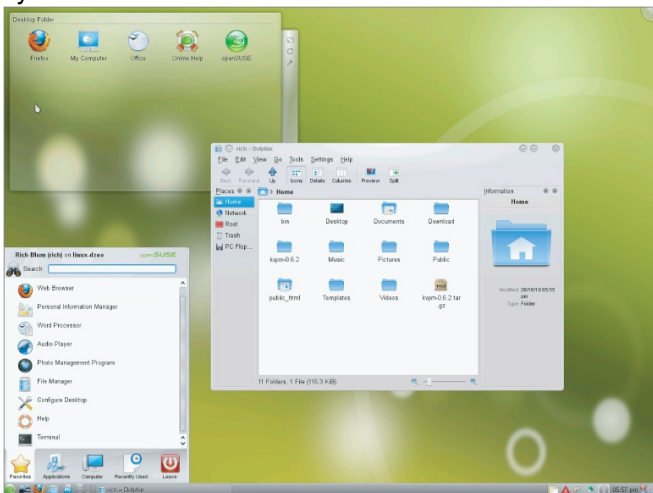
The core X Windows software produces a graphical display environment, but nothing else. While this is fine for running individual applications, it is not too useful for day-to-day computer use. There is no desktop environment allowing users to manipulate files or launch programs. To do that, you need a desktop environment on top of the X Windows system software.

The KDE Desktop

The K Desktop Environment (KDE) was first released in 1996 as an open source project to produce a graphical desktop similar to the Microsoft Windows environment. The KDE desktop incorporates all of the features you are

probably familiar with if you are a Windows user. [Figure 1.3](#) shows a sample KDE 4 desktop running in the openSUSE Linux distribution.

Figure 1.3 The KDE 4 desktop on an openSUSE Linux system



The KDE desktop allows you to place both application and file icons in a special area on the desktop. If you single-click an application icon, the Linux system starts the application. If you single-click on a file icon, the KDE desktop attempts to determine what application to start to handle the file.

The bar at the bottom of the desktop is called the Panel. The Panel consists of four parts:

- **The K menu:** Much like the Windows Start menu, the K menu contains links to start installed applications.
- **Program shortcuts:** These are quick links to start applications directly from the Panel.
- **The taskbar:** The taskbar shows icons for applications currently running on the desktop.
- **Applets:** These are small applications that have an icon in the Panel that often can change depending on information from the application.

All of the Panel features are similar to what you would find in Windows. In addition to the desktop features, the KDE project has produced a wide assortment of applications that run in the KDE environment. These applications are shown in [Table 1.3](#). (You may notice the trend of using a capital K in KDE application names.)

[Table 1.3](#) KDE Applications

Application	Description
amaroK	Audio file player
digiKam	Digital camera software
dolphin	File manager
K3b	CD-burning software
Kaffeine	Video player
Kmail	E-mail client
Koffice	Office applications suite

Konqueror	File and Web browser
Kontakt	Personal information manager
Kopete	Instant messaging client

This is only a partial list of applications produced by the KDE project. There are lots more applications that are included with the KDE desktop.

The GNOME Desktop

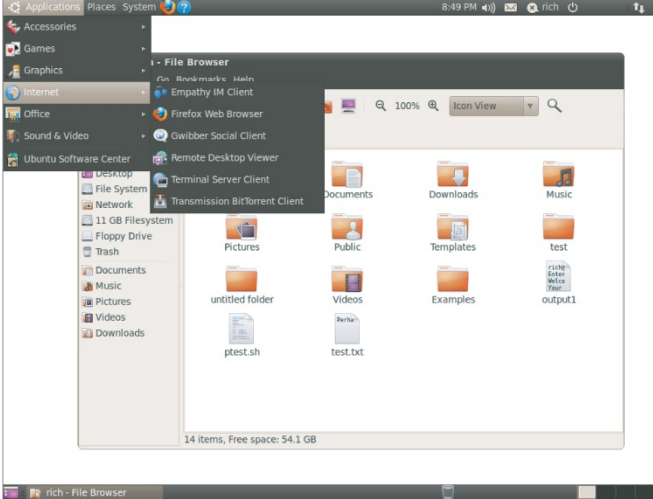
The GNU Network Object Model Environment (GNOME) is another popular Linux desktop environment. First released in 1999, GNOME has become the default desktop environment for many Linux distributions (the most popular being Red Hat Linux).

While GNOME chose to depart from the standard Microsoft Windows look-and-feel, it incorporates many features that most Windows users are comfortable with:

- A desktop area for icons
- Two panel areas
- Drag-and-drop capabilities

[Figure 1.4](#) shows the standard GNOME desktop used in the Ubuntu Linux distribution.

[Figure 1.4](#) A GNOME desktop on an Ubuntu Linux system



Not to be outdone by KDE, the GNOME developers have also produced a host of graphical applications that integrate with the GNOME desktop. These are shown in [Table 1.4](#).

As you can see, there are also quite a few applications available for the GNOME desktop. Besides all of these applications, most Linux distributions that use the GNOME desktop also incorporate the KDE libraries, allowing you to run KDE applications on your GNOME desktop.

Table 1.4 GNOME Applications

Application	Description
-------------	-------------

epiphany	Web browser
evince	Document viewer
gcalc-tool	Calculator
gedit	GNOME text editor
gnome-panel	Desktop panel for launching applications
gnome-nettool	Network diagnostics tool
gnome-terminal	Terminal emulator
nautilus	Graphical file manager
nautilus-cd-burner	CD-burning tool
sound juicer	Audio CD-ripping tool
tomboy	Note-taking software
totem	Multimedia player

Other Desktops

The downside to a graphical desktop environment is that they require a fair amount of system resources to operate properly. In the early days of Linux, a hallmark and selling feature of Linux was its ability to operate on older, less powerful PCs that the newer Microsoft desktop products couldn't run on. However, with the popularity of KDE and GNOME desktops, this has changed, as it takes just as much memory to run a KDE or GNOME desktop as the latest Microsoft desktop environment.

If you have an older PC, don't be discouraged. The Linux developers have banded together to take Linux back to its roots. They've created several low-memory-oriented

graphical desktop applications that provide basic features that run perfectly fine on older PCs.

While these graphical desktops don't have a plethora of applications designed around them, they still run many basic graphical applications that support features such as word processing, spreadsheets, databases, drawing, and, of course, multimedia support.

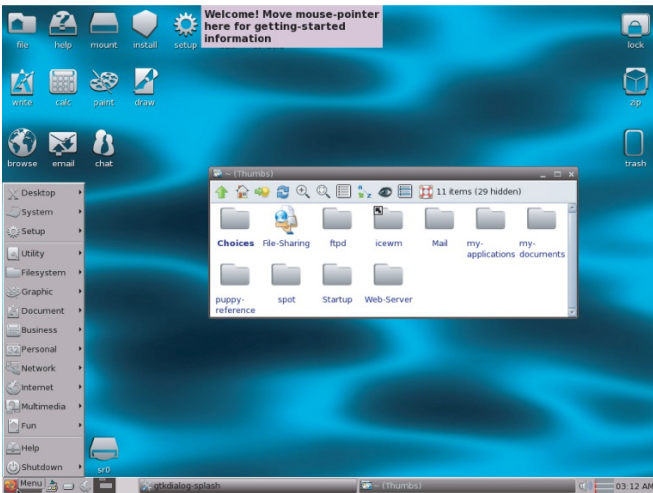
[Table 1.5](#) shows some of the smaller Linux graphical desktop environments that can be used on lower-powered PCs and laptops.

[Table 1.5](#) Other Linux Graphical Desktops

Desktop	Description
fluxbox	A bare-bones desktop that doesn't include a Panel, only a pop-up menu to launch applications
xfce	A desktop that's similar to the KDE desktop, but with less graphics for low-memory environments
JWM	Joe's Window Manager, a very lightweight desktop ideal for low-memory and low-disk space environments
fwm	Supports some advanced desktop features such as virtual desktops and Panels, but runs in low-memory environments
fwm95	Derived from fwm, but made to look like a Windows 95 desktop

These graphical desktop environments are not as fancy as the KDE and GNOME desktops, but they provide basic graphical functionality just fine. [Figure 1.5](#) shows what the fluxbox desktop used in the Puppy Linux antiX distribution looks like.

Figure 1.5 The JWM desktop as seen in the Puppy Linux distribution



If you are using an older PC, try a Linux distribution that uses one of these desktops and see what happens. You may be pleasantly surprised.

Linux Distributions

Now that you have seen the four main components required for a complete Linux system, you may be wondering how you are going to get them all put together to make a Linux

system. Fortunately, there are people who have already done that for you.

A complete Linux system package is called a *distribution*. There are lots of different Linux distributions available to meet just about any computing requirement you could have. Most distributions are customized for a specific user group, such as business users, multimedia enthusiasts, software developers, or average home users. Each customized distribution includes the software packages required to support specialized functions, such as audio- and video-editing software for multimedia enthusiasts, or compilers and integrated development environments (IDEs) for software developers.

The different Linux distributions are often divided into three categories:

- Full core Linux distributions
- Specialized distributions
- LiveCD test distributions

The following sections describe these different types of Linux distributions, and show some examples of Linux distributions in each category.

Core Linux Distributions

A core Linux distribution contains a kernel, one or more graphical desktop environments, and just about every Linux application that is available, precompiled for the kernel. It provides one-stop shopping for a complete Linux

installation. [Table 1.6](#) shows some of the more popular core Linux distributions.

Table 1.6 Core Linux Distributions

Distribution	Description
Slackware	One of the original Linux distribution sets, popular with Linux geeks
Red Hat	A commercial business distribution used mainly for Internet servers
Fedora	A spin-off from Red Hat but designed for home use
Gentoo	A distribution designed for advanced Linux users, containing only Linux source code
Mandriva	Designed mainly for home use (previously called Mandrake)
openSuSe	Different distributions for business and home use
Debian	Popular with Linux experts and commercial Linux products

In the early days of Linux, a distribution was released as a set of floppy disks. You had to download groups of files and then copy them onto disks. It would usually take 20 or more disks to make an entire distribution! Needless to say, this was a painful experience.

Nowadays, with home computers commonly having CD and DVD players built in, Linux distributions are released as either a CD set or a single DVD. This makes installing Linux much easier.

However, beginners still often run into problems when they install one of the core Linux distributions. To cover just about any situation in which someone might want to use Linux, a single distribution has to include lots of application

software. They include everything from high-end Internet database servers to common games. Because of the quantity of applications available for Linux, a complete distribution often takes four or more CDs.

While having lots of options available in a distribution is great for Linux geeks, it can become a nightmare for beginning Linux users. Most distributions ask a series of questions during the installation process to determine which applications to load by default, what hardware is connected to the PC, and how to configure the hardware. Beginners often find these questions confusing. As a result, they often either load way too many programs on their computer or don't load enough and later discover that their computer won't do what they want it to.

Fortunately for beginners, there's a much simpler way to install Linux.

Specialized Linux Distributions

A new subgroup of Linux distributions has started to appear. These are typically based on one of the main distributions but contain only a subset of applications that would make sense for a specific area of use.

In addition to providing specialized software (such as only office products for business users), customized Linux distributions also attempt to help beginning Linux users by autodetecting and autoconfiguring common hardware devices. This makes installing Linux a much more enjoyable process.

[Table 1.7](#) shows some of the specialized Linux distributions available and what they specialize in.

[Table 1.7](#) Specialized Linux Distributions

Distribution	Description
Xandros	A commercial Linux package configured for beginners
SimplyMEPIS	A free distribution for home use
Ubuntu	A free distribution for school and home use
PCLinuxOS	A free distribution for home and office use
Mint	A free distribution for home entertainment use
dyne:bolic	A free distribution designed for audio and MIDI applications
Puppy Linux	A free small distribution that runs well on older PCs

That's just a small sampling of specialized Linux distributions. There are literally hundreds of specialized Linux distributions, and more are popping up all the time on the Internet. No matter what your specialty, you'll probably find a Linux distribution made for you.

Many of the specialized Linux distributions are based on the Debian Linux distribution. They use the same installation files as Debian but package only a small fraction of a full-blown Debian system.

The Linux LiveCD

A relatively new phenomenon in the Linux world is the bootable Linux CD distribution. This lets you see what a Linux system is like without actually installing it. Most

modern PCs can boot from a CD instead of the standard hard drive. To take advantage of this, some Linux distributions create a bootable CD that contains a sample Linux system (called a *Linux LiveCD*). Because of the limitations of the single CD size, the sample can't contain a complete Linux system, but you'd be surprised at all the software they can cram in there. The result is that you can boot your PC from the CD and run a Linux distribution without having to install anything on your hard drive!

This is an excellent way to test various Linux distributions without having to mess with your PC. Just pop in a CD and boot! All of the Linux software will run directly off the CD. There are lots of Linux LiveCDs that you can download from the Internet and burn onto a CD to test drive.

[Table 1.8](#) shows some popular Linux LiveCDs that are available.

Table 1.8 Linux LiveCD Distributions

Distribution	Description
Knoppix	A German Linux, the first Linux LiveCD developed
SimplyMEPIS	Designed for beginning home Linux users
PCLinuxOS	Full-blown Linux distribution on a LiveCD
Ubuntu	A worldwide Linux project, designed for many languages
Slax	Alive Linux CD based on Slackware Linux
Puppy Linux	A full-featured Linux designed for older PCs

You may notice a familiarity in this table. Many specialized Linux distributions also have a Linux LiveCD version. Some Linux LiveCD distributions, such as Ubuntu,

allow you to install the Linux distribution directly from the LiveCD. This enables you to boot with the CD, test drive the Linux distribution, and then if you like it, install it on your hard drive. This feature is extremely handy and user-friendly.

As with all good things, Linux LiveCDs have a few drawbacks. Because you access everything from the CD, applications run more slowly, especially if you're using older, slower computers and CD drives. Also, because you can't write to the CD, any changes you make to the Linux system will be gone the next time you reboot.

But there are advances being made in the Linux LiveCD world that help to solve some of these problems. These advances include the ability to:

- Copy Linux system files from the CD to memory
- Copy system files to a file on the hard drive
- Store system settings on a USB memory stick
- Store user settings on a USB memory stick

Some Linux LiveCDs, such as Puppy Linux, are designed with a minimum number of Linux system files. The LiveCD boot scripts copies them directly into memory when the CD boots. This allows you to remove the CD from the computer as soon as Linux boots. Not only does this make your applications run much faster (because applications run faster from memory), but it also gives you a free CD tray to use for ripping audio CDs or playing video DVDs from the software included in Puppy Linux.

Other Linux LiveCDs use an alternative method that

allows you to remove the CD from the tray after booting. It involves copying the core Linux files onto the Windows hard drive as a single file. After the CD boots, it looks for that file and reads the system files from it. The dyne:bolic Linux LiveCD uses this technique, which is called docking. Of course, you must copy the system file to your hard drive before you can boot from the CD.

A very popular technique for storing data from a live Linux CD session is to use a common USB memory stick (also called a flash drive or a thumb drive). Just about every Linux LiveCD can recognize a plugged-in USB memory stick (even if the stick is formatted for Windows) and read and write files to and from it. This allows you to boot a Linux LiveCD, use the Linux applications to create files, store those files on your memory stick, and then access them from your Windows applications later (or from a different computer). How cool is that?

Summary

This chapter discussed the Linux system, and the basics of how it works. The Linux kernel is the core of the system, controlling how memory, programs, and hardware all interact with one another. The GNU utilities are also an important piece in the Linux system. The Linux shell, which is the main focus of this book, is part of the GNU core utilities. The chapter also discussed the final piece of a Linux system, the Linux desktop environment. Things have

changed over the years, and Linux now supports several graphical desktop environments.

The chapter also discussed the various Linux distributions. A Linux distribution bundles the various parts of a Linux system into a simple package that you can easily install on your PC. The Linux distribution world consists of full-blown Linux distributions that include just about every application imaginable, as well as specialized Linux distributions that only include applications focused on a special function. The Linux LiveCD craze has created another group of Linux distributions that allow you to easily test drive Linux without even having to install it on your hard drive.

In the next chapter, you look at what you need to start your command line and shell scripting experience. You'll see what you need to do to get to the Linux shell utility from your fancy graphical desktop environment. These days that's not always an easy thing.

Chapter 2

Getting to the Shell

In This Chapter

- Terminal emulation
- The terminfo database
- The Linux console
- The xterm terminal
- The Konsole terminal
- The GNOME terminal

In the old days of Linux, all that was available to work with was the shell. System administrators, programmers, and system users all sat at the Linux console terminal entering text commands and viewing text output. These days, with our fancy graphical desktop environments, it's getting harder just to find a shell prompt on the system to work from. This chapter discusses what is required to provide a command line environment, and then walks you through the terminal emulation packages you may run into in the various Linux distributions.

Terminal Emulation

Back before the days of graphical desktops, the only way to interact with a Unix system was through a text *command line interface* (CLI) provided by the shell. The CLI allowed text input only, and could only display text and rudimentary graphics output.

Because of this restriction, output devices did not have to be very fancy. Often a simple dumb terminal was all that was required to interact with the Unix system. A dumb terminal was usually nothing more than a monitor and keyboard (although later on in life they started getting fancier by utilizing a mouse) connected to the Unix system via a communication cable (usually a multi-wire serial cable). This simple combination provided an easy way to enter text data into the Unix system and view text results.

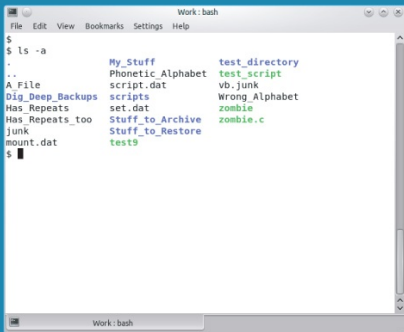
As you well know, things are significantly different in today's Linux environment. Just about every Linux distribution uses some type of graphical desktop environment. However, to access the shell you still need a text display to interact with a CLI. The problem now is getting to one. With all of the new graphical Linux desktop features, sometimes finding a way to get a CLI in a Linux distribution is not an easy task.

One way to get to a CLI is to take the Linux system out of graphical desktop mode and place it in text mode. This provides nothing more than a simple shell CLI on the monitor, just like the days before graphical desktops. This

mode is called the *Linux console* because it emulates the old days of a hard-wired console terminal and is a direct interface to the Linux system.

The alternative to being in the Linux console is to use a *terminal emulation package* from within the graphical Linux desktop environment. A terminal emulation package simulates working on a dumb terminal, all within a graphical window on the desktop. [Figure 2.1](#) shows an example of a terminal emulator running in a graphical Linux desktop environment.

[Figure 2.1](#) A simple terminal emulator running on a Linux desktop



```
Work: bash
File Edit View Bookmarks Settings Help
$ ls -a
.          My_Stuff          test_directory
..         Phonetic_Alphabet test_script
A File     script.dat         vb.junk
Dig_Deep_Backups scripts           Wrong_Alphabet
Has_Repeats set.dat           zombie
Has_Repeats_too Stuff_to_Archive  zombie.c
junk       Stuff_to_Restore
mount.dat  test9
$
```

Each terminal emulation package has the ability to emulate one or more specific types of dumb terminals. If you're going to work with the shell in Linux, unfortunately you'll need to know a little bit about terminal emulation.

Knowing the core features of the old dumb terminals will help you decide which emulation type to select when you're using a graphical terminal emulator, and use all of the available features to their full capabilities. The main features used in the dumb terminal can be broken down into two areas: the graphics capabilities and the keyboard. This section describes these features and discusses how they relate to the different types of terminal emulators.

Graphics Capabilities

The most important part of terminal emulation is how it displays information on the monitor. When you hear the phrase “text mode,” the last thing you'd think to worry about is graphics. However, even the most rudimentary dumb terminals supported some method of screen manipulation (such as clearing the screen and displaying text at a specific location on the screen).

This section describes the graphics features that make each of the different terminal types unique, and what to look for in the terminal emulation packages.

Character Sets

All terminals must display characters on the screen (otherwise, text mode would be pretty useless). The trick is in what characters to display, and what codes the Linux system needs to send to display them. A *character set* is a set of binary commands that the Linux system sends to a monitor to display characters. There are several character sets that are supported by various terminal emulation packages:

- **ASCII:** The American Standard Code for Information Interchange. This character set contains the English characters stored using a 7-bit code, and consists of 128 English letters (both upper and lower case), numbers, and special symbols. This character set

was adopted by the American National Standards Institute (ANSI) as US-ASCII. You will often see it referred to in terminal emulators as the ANSI character set.

- **ISO-8859-1 (commonly called Latin-1):** An extension of the ASCII character set developed by the International Organization for Standardization (ISO). It uses an 8-bit code to support the standard ASCII characters as well as special foreign language characters for most Western European languages. The Latin-1 character set is popular in multinational terminal emulation packages.
- **ISO-8859-2:** ISO character set that supports Eastern European language characters.
- **ISO-8859-6:** ISO character set that supports Arabic language characters.
- **ISO-8859-7:** ISO character set that supports Greek language characters.
- **ISO-8859-8:** ISO character set that supports Hebrew language characters.
- **ISO-10646 (commonly called Unicode):** ISO 2-byte character set that contains codes for most English and non-English languages. This single character set contains all of the codes defined in all of the ISO-8859-x series of character sets. The Unicode character set is quickly becoming popular among open source applications.

By far the most common character set in use today in

English-speaking countries is the Latin-1 character set. The Unicode character set is becoming more popular, and may very well one day become the new standard in character sets. Most popular terminal emulators allow you to select which character set to use in the terminal emulation.

Control Codes

In addition to being able to display characters, terminals must have the ability to control special features on the monitor and keyboard, such as the cursor location on the screen. They accomplish this using a system of *control codes*. A control code is a special code not used in the character set, which signals the terminal to perform a special, nonprintable operation.

Common control code functions are the carriage return (return the cursor to the beginning of the line), line feed (put the cursor on the next horizontal row), horizontal tab (shift the cursor over a preset number of spaces), arrow keys (up, down, left, and right), and the page up/page down keys. While these codes mainly emulate features that control where the cursor is placed on the monitor, there are also several other codes, such as clearing the entire screen, and even a bell ring (emulating the old typewriter end-of-carriage bell).

Control codes were also used in controlling the communication features of dumb terminals. Dumb terminals were connected to the computer system via some type of communication channel, often a serial

communication cable. Sometimes data needed to be controlled on the communication channel, so developers devised special control codes just for data communication purposes. While these codes aren't necessarily required in modern terminal emulators, most support these codes to maintain compatibility. The most common codes in this category are the XON and XOFF codes, which start and stop data transmission to the terminal, respectively.

Block Mode Graphics

As dumb terminals became more popular, manufacturers started experimenting with rudimentary graphics capabilities. By far the most popular type of “graphical” dumb terminal used in the Unix world was the Digital Equipment Corporation (DEC) VT series of terminals. The turning point for dumb terminals came with the release of the DEC VT100 in 1978. The DEC VT100 terminal was the first terminal to support the complete ANSI character set, including block mode graphic characters.

The ANSI character set contains codes that not only allowed monitors to display text but also rudimentary graphics symbols, such as boxes, lines, and blocks. By far one of the most popular dumb terminals used in Unix operations during the 1980s was the DEC VT102, an upgraded version of the VT100. Most modern terminal emulation programs still emulate the operation of the VT102 display, supporting all of the ANSI codes for creating block mode graphics.

Vector Graphics

The Tektronix company produced a popular series of terminals that used a display method called vector graphics. Vector graphics deviated from the DEC method of block mode graphics by making all screen images (including characters) a series of line segments (vectors). The Tektronix 4010 terminal was the most popular graphical dumb terminal produced. Many terminal emulation packages still emulate its capabilities.

The 4010 terminal displays images by drawing a series of vectors using an electron beam, much like drawing with a pencil. Because vector graphics doesn't use dots to create lines, it has the ability to draw geometric shapes using higher precision than most dot-oriented graphics terminals. This was a popular feature among mathematicians and scientists.

Modern terminal emulators use software to emulate the vector graphics drawing capabilities of the Tektronix 4010 terminals. This is still a popular feature for people who need precise graphical drawings, or those who still run applications that used the vector graphics routines to draw complicated charts and diagrams.

Display Buffering

A key to graphics displays is the ability of the terminal to buffer data. Buffering data requires having additional

internal memory within the terminal itself to store characters not currently being displayed on the monitor.

The DEC VT series of terminals utilized two types of data buffering:

- Buffering data as it scrolled off of the main display window (called a history)
- Buffering a completely separate display window (called an alternate screen)

The first type of buffering is known as a *scroll region*. The scroll region is the amount of memory the terminal has that enables it to “remember” data as it scrolls off of the screen. A standard DEC VT102 terminal contained a viewing area for 25 lines of characters. As the terminal displays a new line of characters, the previous line is scrolled upward. When the terminal reaches the bottom line of the display, the next line causes the top line to scroll off the display.

The internal memory in the VT102 terminal allowed it to save the last 64 lines that had scrolled off of the display. Users had the ability to lock the current screen display and use arrow keys to scroll backward through the previous lines that had “scrolled off” of the display. Terminal emulation packages allow you to use either a side scrollbar or a mouse scroll button to scroll through the saved data without having to lock the display. Of course, for full emulation compatibility, most terminal emulation packages also allow you to lock the display and use arrow and page up/page down to scroll through the saved data.

The second type of buffering is known as an *alternative*

screen. Normally, the terminal writes data directly to the normal display area on the monitor. A method was developed to crudely implement animation by using two screen areas to store data. Control codes were used to signal the terminal to write data to the alternative screen instead of the current display screen. That data was held in memory. Another control code would signal the terminal to switch the monitor display between the normal screen data and the data contained in the alternative screen almost instantaneously. By storing successive data pages in the alternative screen area, then displaying it, you could crudely simulate moving graphics.

Terminals that emulate the VT series of terminals have the ability to support both the scroll region and the alternative screen buffering methods.

Color

Even back in the black-and-white (or green) dumb terminal days, programmers were experimenting with different ways to present data. Most terminals supported special control codes to produce the following types of special text:

- Bold characters
- Underline characters
- Reverse video (black characters on white background)
- Blinking
- Combinations of all of the above features

Back in the old days, if you wanted to get someone's attention, you used bold, blinking, reverse video text. Now there's something that could hurt your eyes!

As color terminals became available, programmers added special control codes to display text in various colors and shades. The ANSI character set includes control codes for specifying specific colors for both foreground text and the background color displayed on the monitor. Most terminal emulators support the ANSI color control codes.

The Keyboard

There is more to a terminal than just how the monitor operates. If you have ever worked with different types of dumb terminals, you have seen that they often contain different keys on the keyboard. Trying to emulate specific keys on a specific dumb terminal has proven to be a difficult task for terminal emulation packages.

It was impossible for the creators of the PC keyboard to include keys for every possible type of special key found in dumb terminals. Some PC manufacturers experimented with including special keys for special functions, but eventually the PC keyboard keys became somewhat standardized.

For a terminal emulation package to completely emulate a specific type of dumb terminal, it must remap any dumb terminal keys that don't appear on the PC keyboard. This remapping feature can often become confusing, especially when different systems use different control codes for the

same key.

Some common special keys you'll see in terminal emulation packages are:

- **Break:** Sends a stream of zeroes to the host. This is often used to interrupt the currently executing program in the shell.
- **Scroll Lock:** Also called “no scroll,” this stops the output on the display. Some terminals included memory to hold the contents of the display so the user could scroll backward through previously viewed information while the scroll lock was enabled.
- **Repeat:** When held down with another key, this caused the terminal to repeatedly send the other key's value to the host.
- **Return:** Commonly used to send a carriage return character to the host. Most often used to signify the end of a command for the host to process (now called Enter on PC keyboards).
- **Delete:** While basically a simple feature, the Delete key causes grief for terminal emulation packages. Some terminals delete the character at the current cursor location, while others delete the preceding character. To resolve this dilemma, PC keyboards include two delete keys, Backspace and Delete.
- **Arrow keys:** Commonly used to position the cursor at a specific place—for example, when scrolling through a listing.

- **Function keys:** A combination of specialty keys that can be assigned unique values in programs similar to the PC F1 through F12 keys. The DEC VT series of terminals actually had two sets of function keys, F1 through F20, and PF1 through PF4.

Keyboard emulation is a crucial element in a terminal emulation package. Unfortunately, often applications are written requiring users to hit specific keys for specific functions. I've seen many a communications package that used the old DEC PF1 through PF4 keys, which are often hard to find on a terminal emulation keyboard.

The terminfo Database

Now that you know about terminal emulation packages that can emulate different types of terminals, you need a way for the Linux system to know exactly what terminal you're emulating. The Linux system needs to know what control codes to use when communicating with the terminal emulator. This is done by using an environment variable (see Chapter 5) and a special set of files collectively called the *terminfo database*.

The terminfo database is a set of files that identify the characteristics of various terminals that can be used on the Linux system. The Linux system stores the terminfo data for each terminal type as a separate file in the terminfo database directory. The location of this directory often varies from distribution to distribution. Some common

locations are `/usr/share/terminfo`, `/etc/terminfo`, and `/lib/terminfo`.

To help with organization (often there are lots of different terminfo files), you will see that the terminfo database directory contains directories for different letters of the alphabet. The individual files for specific terminals are stored under the appropriate letter directory for their terminal name. For example, under `/usr/share/terminfo/v` are the VT terminal emulators.

An individual terminfo file is a binary file that is the result of compiling a text file. This text file contains code words that define screen functions, associated with the control code required to implement the function on the terminal.

Since the terminfo database files are binary, you cannot see the codes within these files. However, you can use the `infocmp` command to convert the binary entries into text. An example of using this command is:

```
$ infocmp vt100
# Reconstructed via infocmp from file: /lib/terminfo/v/vt100
vt100|vt100-am|dec vt100 (w/advanced video),
    am, msgr, xenl, xon,
    cols#80, it#8, lines#24, vt#3,
    acsc="aaffggjjkkllmmnnoppqrrssttuuvvwwxxyyzz{|}|}~",
    bel=^G, blink=\E[5m$<2>, bold=\E[1m$<2>,
    clear=\E[H\E[J$<50>, cr=^M, csr=\E[%p1$d;%p2$dr,
    cub=\E[%p1$dD, cubl=^H, cud=\E[%p1$dB, cudl=^J,
    cuf=\E[%p1$dC, cuf1=\E[C$<2>,
    cup=\E[%i%p1$d;%p2$dH$<5>, cuu=\E[%p1$dA,
    cuul=\E[A$<2>, ed=\E[J$<50>, el=\E[K$<3>, ell=\E[1K$<3>,
    enacs=\E(B\E)O, home=\E[H, ht=^I, hts=\EH, ind=^J, kal=\EOq,
    ka3=\EOs, kb2=\EOr, kbs=^H, kcl=\EOp, kc3=\EOn, kcu1=\EOD,
    kcud1=\EOB, kcufl=\EOC, kcuul=\EOA, kent=\EOM, kf0=\EOy,
    kfl=\EOP, kf10=\EOx, kf2=\EOQ, kf3=\EOR, kf4=\EOS, kf5=\EOt,
    kf6=\EOu, kf7=\EOv, kf8=\EOl, kf9=\EOw, rc=\E8,
    rev=\E[7m$<2>, ri=\EM$<5>, rmacs=^O, rmam=\E[?71,
    rmkx=\E[?11\E>, rmso=\E[m$<2>, rmul=\E[m$<2>,
```

```
rs2=\E>\E[?1\E[?41\E[?51\E[?7h\E[?8h,
sc=\E7,
sgr0=\E[m\017$<2>, smacs='N, smam=\E[?7h, smkx=\E[?1h\E=,
smso=\E[7m$<2>, smul=\E[4m$<2>, tbc=\E[3g,
```

\$

The terminfo entry defines the terminal name (in this case vt100), along with any alias names that can be associated with the terminal name. Notice that the first line shows the location of the terminfo file the values were extracted from.

Following that, the `infocmp` command lists the capabilities of the terminal definition, along with the control codes used to emulate the individual capabilities. Some capabilities are either enabled or disabled. If the capability appears in the list, it's enabled by the terminal definition (such as the `am`, auto-right-margin, feature). Other capabilities must define a specific control code sequence to perform the task (such as clearing the monitor display). [Table 2.1](#) shows a list of some of the capabilities you see in the vt100 terminfo definition file listed.

Table 2.1 Terminfo Capability Codes

Code	Description
<code>am</code>	Set right-side auto-margin
<code>msgsr</code>	Safe to move cursor in standout mode
<code>xenl</code>	Newline characters ignored after 80 columns
<code>xon</code>	Terminal uses XON/XOFF characters for flow control
<code>cols#80</code>	80 columns in a line
<code>it#8</code>	Tab character set to eight spaces
<code>lines#24</code>	24 lines on a screen
<code>vt#3</code>	Virtual terminal number 3

bel	Control code to use to emulate the bell
blink	Control code used to produce blinking text
bold	Control code used to produce bold text
clear	Control code used to clear the screen
cr	Control code used to enter a carriage return
csr	Control code used to change scroll region
cub	Move one character to the left without erasing
cubl	Move cursor back one space
cud	Move cursor down one line
cudl	Control code to move cursor down one line
cuf	Move one character to the right without erasing
cubl	Control code to move the cursor right one space without erasing
cup	Control code to move to row one, column two on the display
cuu	Move cursor up one line
cuul	Control code to move cursor up one line
ed	Clear to the end of the screen
el	Clear to the end of the line
ell	Clear to the beginning of the line.
enacs	Enable the alternate character set
home	Control code to move cursor to the home position—row one, column two (same as cup)
ht	Tab character
hts	Set tab in every row at current column
ind	Scroll text up
ka1	Upper-left key in keypad
ka3	Upper-right key in keypad

kb2	Center key in keypad
kbs	Backspace key
kc1	Lower-left key in keypad
kc3	Lower-right key in keypad
kcub1	The left arrow key
kcud1	Control code for down arrow key
kcuf1	The right arrow key
kcuu1	The up arrow key
kent	The Enter key
kf0	The F0 function key
kf1	The F1 function key
kf10	The F10 function key
rc	Restore cursor to last saved position
rev	Reverse video mode
ri	Scroll text down
rmacs	End alternate character set
rmam	Turn off automatic margins
rmkx	Exit keyboard transmit mode
rmso	Exit standout mode
rmul	Exit underline mode
rs2	Reset
sc	Save current cursor position
sgr	Define video attributes
sgr0	Turn off all attributes
smacs	Start alternate character set
smam	Turn on automatic margins

smkx	Start keyboard transmit mode
sms0	Begin standout mode
smul	Begin underline mode
tbc	Clear all tab stops

The Linux shell uses the `TERM` environment variable to define which terminal emulation setting in the terminfo database to use for a specific session. When the `TERM` environment variable is set to `vt100`, the shell knows to use the control codes associated with the `vt100` terminfo database entry for sending control codes to the terminal emulator. To see the `TERM` environment variable, you can just echo it from the CLI:

```
$ echo $TERM
xterm
$
```

This example shows that the current terminal type is set to the `xterm` entry in the terminfo database.

The Linux Console

In the early days of Linux, when you booted up your system you would see a login prompt on your monitor, and that's all. As mentioned earlier, this is called the Linux console. It was the only place you could enter commands for the system.

With modern Linux systems, when the Linux system starts it automatically creates several *virtual consoles*. A virtual console is a terminal session that runs in memory on the Linux system. Instead of having several dumb terminals

connected to the PC, most Linux distributions start seven (or sometimes even more) virtual consoles that you can access from the single PC keyboard and monitor.

In most Linux distributions, you can access the virtual consoles using a simple keystroke combination. Usually you must hold down the Ctl+Alt key combination, and then press a function key (F1 through F8) for the virtual console you want to use. Function key F1 produces virtual console 1, key F2 produces virtual console 2, and so on.

Six of the virtual consoles use a full-screen text terminal emulator to display a text login screen, as shown in [Figure 2.2](#).

[Figure 2.2](#) The Linux console login screen

```
Ubuntu 10.10 user-desktop tty2
user-desktop login: user
Password:
Last login: Mon Dec 20 09:03:13 EST 2010 on tty2
Linux user-desktop 2.6.35-24-generic #42-Ubuntu SMP Thu Dec 2 01:41:57 UTC 2010 i686 GNU/Linux
Ubuntu 10.10

Welcome to Ubuntu!
 * Documentation: https://help.ubuntu.com/

No mail.
$
```

After logging in with your user ID and password, you are taken to the Linux bash shell CLI. In the Linux console, you do not have the ability to run any graphical programs. You can only use text programs to display on the Linux text consoles.

After logging in to a virtual console, you can keep it active and switch to another virtual console without losing your active session. You can switch between all of the virtual consoles, with multiple active sessions running.

The first or the last two virtual consoles are normally reserved for X Window graphical desktops. Some

distributions only assign one so you may have to test all three Ctl+Alt+F1, Ctl+Alt+F7, and Ctl+Alt+F8 to see which one your particular distribution uses. Most distributions automatically switch to one of the graphical virtual consoles after the boot sequence completes, providing a complete graphical login and desktop experience.

Logging in to a text virtual terminal session and then switching over to a graphical one can get tedious. Fortunately, there's a better way to jump between graphical and text mode on the Linux system: terminal emulation packages are a popular way to access the shell CLI from within a graphical desktop session. The following sections describe the most common software packages that provide terminal emulation in a graphical window.

The xterm Terminal

The oldest and most basic of X Window terminal emulation packages is *xterm*. The xterm package has been around since the original days of X Window, and is included by default in most X Window packages.

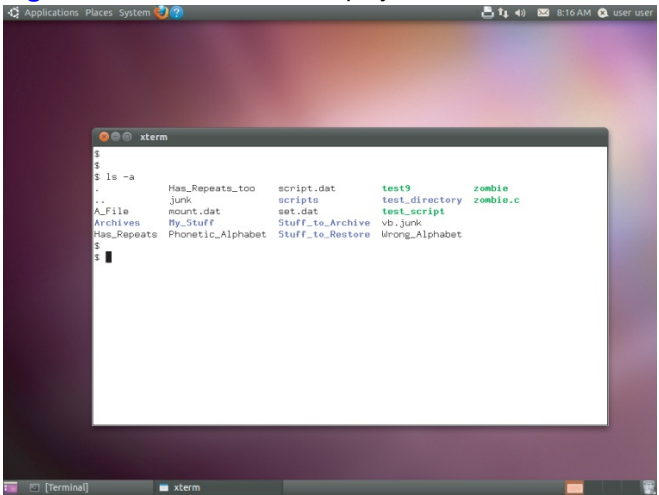
The xterm package provides both a basic VT102/220 terminal emulation CLI and a graphical Tektronix 4014 environment (similar to the 4010 environment). While xterm is a full terminal emulation package, it doesn't require many resources (such as memory) to operate. Because of this, the xterm package is still popular in Linux distributions designed to run on older hardware. Some graphical

desktop environments, such as fluxbox, use it as the default terminal emulation package.

While not offering many fancy features, the xterm package does one thing extremely well, and that is emulate a VT220 terminal. The newer versions of xterm even emulate the VT series of color control codes, allowing you to use color in your scripts.

[Figure 2.3](#) shows what the basic xterm display looks like running on a graphical Linux desktop.

Figure 2.3 The basic xterm display



The xterm package allows you to set individual features

using both command line parameters and a series of four simple graphical menus. The following sections discuss these features and how to change them.

Command Line Parameters

The list of xterm command line parameters is extensive. There are lots of features you can control to customize the terminal emulation features, such as enabling or disabling individual VT emulations.

The xterm command line parameters use the plus (+) and minus (-) signs to signify how a feature is set. A plus sign indicates that the feature should be returned to the default setting. A minus sign indicates that you are setting the feature to a non-default value. [Table 2.2](#) lists some of the more common features that you can set using the command line parameters.

Table 2.2 xterm Command Line Parameters

Parameter	Description
<code>l32</code>	By default, xterm does not allow 132 characters per line mode.
<code>ah</code>	Always highlight the text cursor.
<code>aw</code>	Auto-line-wrap is enabled.
<code>bc</code>	Enables text cursor blinking.
<code>bg color</code>	Specify the color to use for the background.
<code>cm</code>	Disables recognition of ANSI color change control codes.
<code>fb font</code>	Specify the font to use for bold text.

<code>fg color</code>	Specify the color to use for the foreground text.
<code>fn font</code>	Specify the font to use for text.
<code>fw font</code>	Specify the font to use for wide text.
<code>hc color</code>	Specify the color to use for highlighted text.
<code>j</code>	Use jump scrolling, scrolling multiple lines at a time.
<code>l</code>	Enable logging screen data to a log file.
<code>lf filename</code>	Specify the file name to use for screen logging.
<code>mb</code>	Ring a margin bell when the cursor reaches the end of a line.
<code>ms color</code>	Specify the color used for the text cursor.
<code>name name</code>	Specify the name of the application that appears in the title bar.
<code>rv</code>	Enable reverse video by swapping the background and foreground colors.
<code>sb</code>	Use a side scrollbar to allow scrolling of saved scroll data.
<code>t</code>	Start xterm in Tektronix mode.
<code>tb</code>	Specify that xterm should display a toolbar at the top.

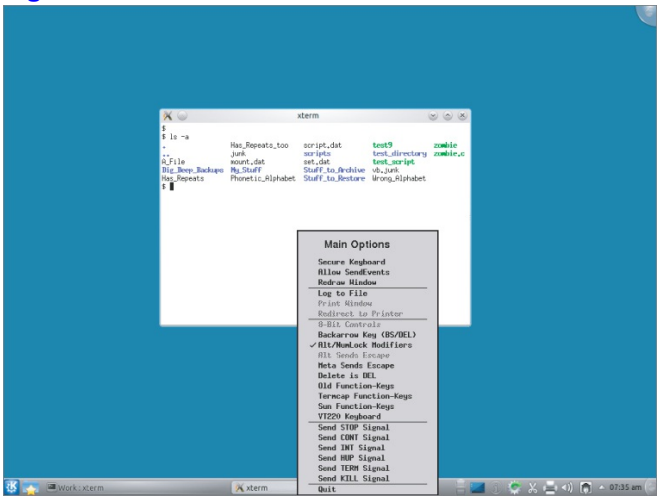
It is important to note that not all implementations of xterm support all of these command line parameters. You can determine which parameters your xterm implements by using the `-help` parameter when you start xterm on your system.

The xterm Main Menu

The main xterm menu contains configuration items that apply to both the VT102 and Tektronix windows. You can access the main menu by holding down the Ctrl key and

clicking the mouse button once (the left button on a right-hand mouse, the right button on a left-hand mouse) while in an xterm session window. [Figure 2.4](#) shows what the xterm main menu looks like.

Figure 2.4 The xterm main menu



There are four sections in the xterm main menu, as described in the following sections.

X Event Commands

The X event commands section contains features that allow

you to manage how xterm interacts with the X Window display.

- **Toolbar:** If the xterm installation supports the toolbar, this entry enables or disables displaying the toolbar in the xterm window (the same as the `tb` command line parameter).
- **Secure Keyboard:** Restricts the keyboard keystrokes to a specific xterm window. This is useful when typing passwords to ensure they don't get hijacked by another window.
- **Allow SendEvents:** Allows X Window events generated by other X Window applications to be accepted by the xterm window.
- **Redraw Window:** Instructs X Window to refresh the xterm window.

Again, all of these features may not be supported by your particular xterm implementation. If they're not supported, they'll appear grayed-out in the menu.

Output Capturing

The xterm package allows you to capture data displayed in the window and either log it to a file or send it to a default printer defined in X Window. The features that appear in this section are:

- **Log to file:** Sends all data displayed in the xterm window to a text file.
- **Print window:** Sends all data displayed in the

current window to the default X Window printer.

- **Redirect to printer:** Sends all data displayed in the xterm window to the default X Window printer as well. This feature must be turned off to stop printing data.

The capturing feature can get messy if you are using graphics characters or control characters (such as colored text) in your display area. All characters sent to the display, including control characters, are stored in the log file or sent to the printer.

The xterm print feature assumes that you define a default printer in the X Window system. If you have no printer defined, the feature will appear grayed out in the menu.

Keyboard Settings

The keyboard settings section contains features that allow you to customize how xterm sends keyboard characters to the host system.

- **8-bit controls:** Sends 8-bit control codes, used in VT220 terminals, rather than 7-bit ASCII control codes.
- **Back arrow key:** Toggles the back arrow key between sending a backspace character or a delete character.
- **Alt/Numlock Modifiers:** Controls whether the Alt or Numlock keys change the PC number pad behavior.
- **Alt Sends Escape:** The Alt key sends an escape control code along with the other key pressed.

- **Meta sends Escape:** Controls whether the function keys send a two-character control code, including the escape control code.
- **Delete is DEL:** The PC Delete key sends a delete character instead of a backspace character.
- **Old Function keys:** The PC functions keys emulate the DEC VT100 function keys.
- **Termcap Function keys:** The PC function keys emulate the Berkley Unix Termcap function keys.
- **Sun Function keys:** The PC function keys emulate the Sun Workstation function keys.
- **VT220 keyboard:** The PC function keys emulate the DEC VT220 function keys.

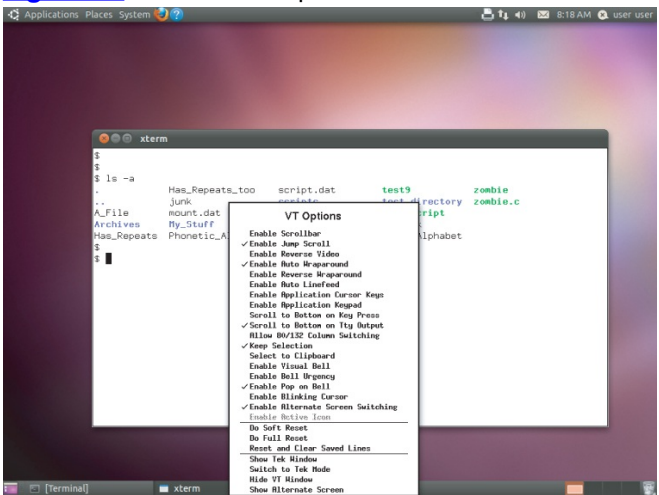
As you can see, setting keyboard preferences often depends on the specific application and/or environment you're working in. There's also a fair amount of personal preference involved as well. Often it's just a matter of what works best for you as to which keyboard settings to make.

The VT Options Menu

The VT options menu sets features xterm uses in the VT102 emulation. You access the VT options menu by holding down the Control key and clicking the second mouse button. Typically, the second mouse button is the middle mouse button. If you're using a two-button mouse, most Linux X Window configurations emulate the middle mouse button when you click both the left and right mouse

buttons together. [Figure 2.5](#) shows what the VT options menu looks like.

Figure 2.5 The xterm VT options menu



As you can see from [Figure 2.5](#), many of the VT features that you can set from the command line parameters can also be set from the VT options menu. This produces quite a large list of available options. The VT options are divided into three sets of commands, described in the following sections.

VT Features

The VT features commands to change the features of how xterm implements the VT102/220 emulation. They include:

- Enable Scrollbar
- Enable Jump Scrollbar
- Enable Reverse Video
- Enable Auto Wraparound
- Enable Reverse Wraparound
- Enable Auto Linefeed
- Enable Application Cursor Keys
- Enable Application Keypad
- Scroll to Bottom on Keypress
- Scroll to Bottom on TTY Output
- Allow 80/132 Column Switching
- Select to Clipboard
- Enable Visual Bell
- Enable Bell Urgency
- Enable Pop on Bell
- Enable Blinking Cursor
- Enable Alternate Screen Switching
- Enable Active Icon

You can enable or disable each of these features by clicking on the feature in the menu. An enabled feature will have a checkmark next to it.

VT Commands

The VT commands section sends a specific reset

command to the xterm emulation window. They include:

- Do Soft Reset
- Do Full Reset
- Reset and Clear Saved Lines

The soft reset sends a control code to reset the screen area. This is convenient if a program sets the scroll region incorrectly. The full reset clears the screen, resets any set tab positions, and resets any terminal mode feature set during the session to the initial state. The Reset and Clear Saved Lines command performs a full reset and also clears out the scroll area history file.

Current Screen Commands

The current screen commands section sends commands to the xterm emulator that affect which screen is the currently active screen.

- **Show Tek Window:** Display the Tektronix terminal window along with the VT100 terminal window.
- **Switch to Tek Window:** Hide the VT100 terminal window and display the Tektronix terminal window.
- **Hide VT Window:** Hide the VT100 terminal window while displaying the Tektronix terminal window.
- **Show Alternate Screen:** Display the data currently stored in the VT100 alternate screen area.

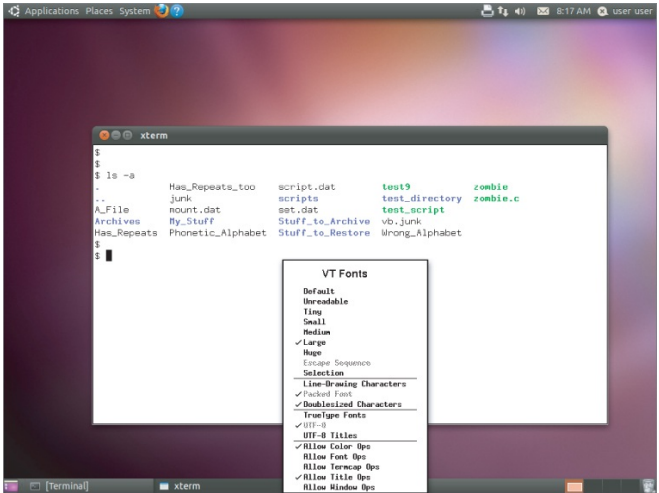
The xterm terminal emulator provides the ability to start in either VT100 terminal mode (by default) or in the Tektronix terminal mode (by using the `t` command line parameter).

After you start in either mode, you can use this menu area to switch to the other mode during your session.

The VT Fonts Menu

The VT fonts menu sets the font style used in the VT100/220 emulation window. You can access this menu by holding the Control key and clicking on mouse button three (the right button on a right-handed mouse, or the left button on a left-handed mouse). [Figure 2.6](#) shows what the VT fonts menu looks like.

Figure 2.6 The xterm VT fonts menu



The VT fonts menu, covered in the following sections, contains three sections of selections.

Set the Font

These menu options set the size of the font used in the xterm window. The available sizes are:

- Default
- Unreadable
- Tiny
- Small
- Medium
- Large
- Huge
- Escape Sequence
- Selection

The default font is the standard-sized font used to display text in the current X Window frame. The unreadable font is pretty much what it says. It shrinks the xterm window down to a size that is not really usable. This is handy, however, when you want to minimize the window on your desktop without completely minimizing it on the system. The large and huge font options produce extremely large font sizes for visually impaired users.

The Escape the Sequence option sets the font to the last font set by the VT100 set font control code. The Selection option allows you to save the current font with a special font

name.

Display the Font

This section of menu options defines the type of characters used to create the text. There are three options available:

- **Line Drawing Characters:** Tells the Linux system to produce ANSI graphical lines instead of using line characters from the chosen font.
- **Packed Font:** Tells the Linux system to use a packed font.
- **Doublesized characters:** Tells the Linux system to scale the set font to double the normal size.

The line drawing characters allow you to determine which types of graphical features to use when drawing in text mode. You can use either characters provided by the selected font source or characters provided by the DEC VT100 control codes.

Specify the Font

This section of the menu provides options for what type of fonts are used to create the characters:

- TrueType Fonts
- UTF-8 Fonts
- UTF-8 Titles

The TrueType fonts are popular in graphical environments. Instead of each character taking the same

amount of space in the line, characters are proportioned by their natural sizes. Thus, the letter *i* takes up less space on the line than the letter *m*. The UTF-8 font allows you to temporarily switch to use the Unicode character set for applications that don't support foreign characters. The Titles option allows the xterm window title to be encoded using UTF-8.

The Konsole Terminal

The KDE Desktop Project has created its own terminal emulation package called *Konsole*. The Konsole package incorporates the basic xterm features, along with more advanced features that we now expect from a Windows application. This section describes the features of the Konsole terminal, and shows how to use them.

Command Line Parameters

Often a Linux distribution provides a method for starting applications directly from the graphical desktop menu system. If your distribution doesn't provide this feature, you can manually start Konsole by using the format:

```
konsole parameters
```

Just like xterm, the Konsole package uses command line parameters to set features in the new sessions. [Table 2.3](#) shows the available Konsole command line parameters.

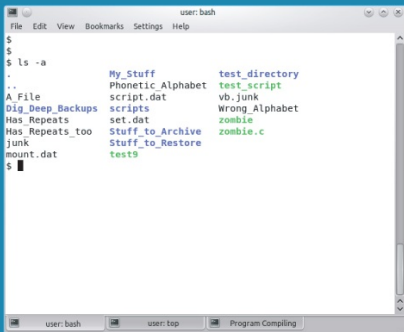
Table 2.3 The Konsole Command Line Parameters

Parameter	Description
-e command	Execute command instead of a shell.
--keytab file	Use the specified key file to define key mappings.
--keytabs	List all of the available keytabs.
--ls	Start the Konsole session with a login screen.
--name name	Set the name that appears in the Konsole title bar.
--noclose	Prevent the Konsole window from closing when the last session has been closed.
--noframe	Start Konsole without a frame.
--nohist	Prevent Konsole from saving scroll history in sessions.
--nomenubar	Start Konsole without the standard menu bar options.
--noresize	Prevent changing the size of the Konsole window area.
--notabbar	Start Konsole without the standard tab area for sessions.
--noxft	Start Konsole without support for aliasing smaller fonts.
--profile file	Start Konsole with settings saved in the specified file.
--profiles	List all of the available Konsole profiles.
--schema name	Start Konsole using the specified schema name or file.
--schemata	List the schemes available in Konsole.
-T title	Set the Konsole window title.
--type type	Start a Konsole session using the specified type.
--types	List all of the available Konsole session types.
--vt_sz CxL	Specify the terminal columns (C) and rows (L).
--workdir dir	Specify the working directory for Konsole to store temporary files.

Tabbed Window Sessions

When you start Konsole, you'll notice that it has a tabbed window, with one tab open to a terminal emulation session. This is the default tabbed window session, and it is normally a standard bash shell CLI. Konsole allows you to have multiple tabs active at the same time. The tabs, placed at either the top or bottom of the window area, allow you to easily switch between sessions. This is a great feature for programmers who need to edit code in one tab, while testing the code in another tabbed window. It's easy to flip back and forth between different active tabs in Konsole. [Figure 2.7](#) shows a Konsole window with three active tabs.

[Figure 2.7](#) The Konsole terminal emulator with three active sessions



```
user: bash
File Edit View Bookmarks Settings Help
$
$
$ ls -a
.          My_Stuff          test_directory
..         Phonetic_Alphabet test_script
A File     script.dat             vb.junk
Dig_Deep_Backups scripts              Wrong_Alphabet
Has_Repeats set.dat              zombie
Has_Repeats_too Stuff_to_Archive     zombie.c
junk        Stuff_to_Restore
mount.dat   test9
$
```

Similar to the xterm terminal emulator, Konsole provides a simple menu by right-clicking in the active tab area. If you right-click in the tab area, a menu appears with the following options:

- **Copy:** Copy the selected text to the clipboard.
- **Paste:** Paste the contents of the clipboard to the selected area.
- **Clear Scrollback & Reset:** Clears all the text out of the current tab and resets the terminal.
- **Open File Manager:** Opens the KDE default file manager, Dolphin, at the present working directory.

- **Change Profile:** Changes the profile for the current tab.
- **Edit Current Profile:** Edits the current tab's profile.
- **Show Menu Bar:** Toggles on/off the menu bar display.
- **Character Encoding:** Selects the character set used to send and display characters.
- **Close Tab:** Terminate the tabbed window session. If it is the last tab in the Konsole window, Konsole will close.

Konsole also provides another quick way to access the new tab menu—by holding down the Ctrl key and right-clicking in the tab area.

After a tabbed window has been modified, you can keep the modifications to use in the future using a profile.

Profiles

Konsole delivers a powerful method, called profiles, for saving and reusing a tabbed session's settings. When you start Konsole for the first time, the tab session's settings are pulled from the default profile, Shell. These settings include items such as what shell to use, color schemes, and so on. Once you have modified your current tab session, you can save those modifications as a new profile. This feature allows multiple tab setups, such as a tab session that uses a different shell than bash.

Profiles can also be used to automate mundane tasks,

such as logging into another system. You can define many profiles and use different ones in each open tab session. To create a new profile, use the Edit Current Profile setting, described previously in the simple Konsole menu. To switch the current profile to a different profile, use the simple menu option Change Profile. These options are also provided on the menu bar.

By default, Konsole uses a menu bar to provide additional functionality so that you can modify and save your Konsole tabs and profiles.

The Menu Bar

The default Konsole setup uses a menu bar for you to easily view and change options and features in your tabs. The menu bar consists of six items, as described in the following sections.

File

The File menu bar item provides a location for starting a new tab in the current window or in a new window. It contains the following entries:

- **New Tab:** Start a new Konsole tab within the current terminal window using the default profile, Shell.
- **New Window:** Start a new terminal window to hold a new Konsole tab.
- **List of defined profiles:** Switch to a new profile in the current tab session.

- **Open File Manager:** Open the file manager at the present working directory.
- **Close Tab:** Close the current tab.
- **Quit:** Quit the Konsole application.

When you first start Konsole, the only profile listed in the List-of-Defined-Profiles will be Shell. As more profiles are created and saved, their names will appear in the list.

Edit

The Edit menu bar provides options for handling text in the session as well as a few additional options:

- **Copy:** Copies selected text (which was highlighted with the mouse) to the system clipboard.
- **Paste:** Pastes text currently in the system clipboard to the current cursor location. If the text contains newline characters, they will be processed by the shell.
- **Rename Tab:** Changes the current tab name. The following tokens can be used in addition to text:
 - %**#**—Session number
 - %**D**—Current directory (absolute name)
 - %**d**—Current directory (relative name)
 - %**n**—Program name
 - %**u**—User name
 - %**w**—Shell set window title
- **Copy Input To:** Sends typed text in the current tab to

one or more tabs in the current terminal windows.

- **ZModem Upload:** Uploads a file to the system using the ZModem protocol.
- **Clear and Reset:** Sends the control code to reset the terminal emulator, and clears the current session window.

Konsole provides an excellent method for tracking a tab's function. Using the Rename Tab menu option, you can name a tab to match its profile. This helps in tracking which open tab is performing what function.

View

The View menu bar item contains items for controlling the individual sessions in the Konsole window. These selections include:

- **Split View:** Controls the display within the current terminal emulation window. Views can be modified by:
 - **Split View Left/Right:** Splits the current display into two identical screens, side by side
 - **Split View Top/Bottom:** Splits the current display into two identical screens, one on top of the other
 - **Close Active:** Merges the current split terminal window back into a single window
 - **Close Others:** Merges non-current split terminal windows back into a single window

- **Expand View:** Adjusts the active split of a terminal window to take up more of the display window
- **Shrink View:** Adjusts the active split of a terminal window to take up less display window
- **Detach View:** Remove the current tab from the Konsole window, and start a new Konsole window using the current tab. This is available only when more than one active tab is open.
- **Show Menu Bar:** Toggles on/off the display of the menu bar.
- **Full Screen Mode:** Toggles on/off the terminal window filling the entire monitor display area.
- **Monitor for Silence:** Toggles on/off a special icon appearance when no new text appears in the tab for 10 seconds. This allows you to switch to another tab while waiting for output from an application to stop, such as when compiling a large application.
- **Monitor for Activity:** Toggles on/off a special icon appearance when new text appears in the tab. This allows you to switch to another tab while waiting for output from an application.
- **Character Encoding:** Selects the character set used to send and display characters.
- **Increase Text Size:** Increases the size of the text font.
- **Decrease Text Size:** Decreases the size of the text

font.

The Split View option in Konsole will maintain the current number of open tabs in the split view. For example, if you have three tabs in the terminal window and split the view, each view will have three tabs.

Scrollback

Konsole retains a history area, formally called a *scrollback buffer*, for each tab. The history area contains the output text for lines that scroll out of the viewing area of the terminal emulator. By default, the last 1,000 lines of output in the scrollback buffer are retained. The Scrollback menu offers various options for reviewing this buffer.

- **Search Output:** Opens a dialog box at the bottom of the current tab. The Find dialog box enables Konsole to search for specific text in the scroll buffer. It has options for case, regular expressions, and search direction.
- **Find Next:** Finds the next text match in more recent history of the scrollback buffer.
- **Find Previous:** Finds the next text match in more ancient history of the scrollback buffer.
- **Save Output:** Saves the contents of the scrollback buffer to a text or HTML file.
- **Scrollback Options:** Controls the activity of the scrollback buffer. The modifications available are:
 - **No Scrollback:** Disable the scrollback buffer.

- **Fixed Scrollback:** Set the size (number of lines) of the scrollback buffer. The default is 1,000 lines.
- **Unlimited Scrollback:** Allows an infinite number of lines to be stored in the scrollback buffer.
- **Save to Current Profile:** Save the scrollback buffer options to the current profile settings.
- **Clear Scrollback & Reset:** Removes the contents of the scrollback buffer and resets the terminal window.

You can scroll through the scrollback buffer by using the scrollbar in the viewing area, or by pressing the Shift key and the Up Arrow key to scroll line by line, or the Page Up key to scroll page (24 lines) by page.

Bookmarks

The Bookmarks menu items provide a way to manage bookmarks set in the Konsole window. A bookmark enables you to save your directory location in an active session and then easily return there in either the same session or a new session. Have you ever drilled down several directories deep to find something on the Linux system, exited, and then forgotten how you got there? Bookmarks will solve that problem. When you get to your desired directory location, just add a new bookmark. When you want to return, look at the Bookmarks for your new

bookmark, and it will automatically perform the directory change to the desired location for you. The bookmark entries include:

- **Add Bookmark:** Create a new bookmark at the current directory location.
- **Bookmark Tabs as Folder:** Create a bookmark for all the current terminal window tabs.
- **New Bookmark Folder:** Create a new storage folder for bookmarks.
- **Edit Bookmarks:** Edit existing bookmarks.
- **A list of your bookmarks:** Any bookmarks you have created.

There is no limit to how many bookmarks you can store in Konsole, but having lots of bookmarks can get confusing. By default, they all appear at the same level in the Bookmark area. You can organize them by creating new bookmark folders and moving individual bookmarks to the new folders using the Edit Bookmarks item.

Settings

The Settings menu bar area allows you to customize and manage your profiles as well as add a little more functionality to the current tab session. This area includes:

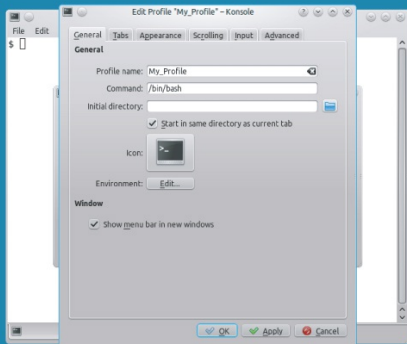
- **Change Profile:** Applies a selected profile for the current tab.
- **Edit Current Profile:** Opens a dialog box where a large variety of profile settings can be changed.

- **Manage Profiles:** Allows a particular profile to be the default profile and enables you to create and delete profiles. Also manages the order in which your profiles appear in the File menu.
- **Configure Shortcuts:** Creates keyboard shortcuts for Konsole commands.
- **Configure Notifications:** Sets actions for specific session events.
- **Configure Konsole:** Create custom Konsole schemas and sessions.

The Configure Notifications area is pretty cool. It allows you to associate five specific events that can occur within a session with six different actions. When one of the events occurs, the defined action (or actions) are taken.

The Edit Current Profile settings is a powerful tool that provides advanced control over profile features. This dialog box provides a way to create and save a variety of profiles for later use. [Figure 2.8](#) shows the main Edit Current Profile dialog box.

[Figure 2.8](#) The Konsole edit current profile dialog box



Within the Edit Current Profile dialog box are six tabbed areas:

- **General:** Allows you to set the profile's name, the icon, the initial filesystem directory, and so on. You can also designate the command to be executed upon the tab opening. This typically points to the bash shell, `/bin/bash`, but can also be regularly used shell commands, such as `top`.
- **Tabs:** Defines the tab's title format and the position of the tab bars.
- **Appearance:** Items such as a tab's color scheme

and font settings are in this window.

- **Scrolling:** Settings included are the size of the scrollbar buffer and the scroll bar's location on the window.
- **Input:** Keybindings, what characters are sent to the terminal emulation when certain keyboard combinations are pressed, can be set here.
- **Advanced:** Allows you to configure several settings in this window. They include terminal features, character encoding, mouse interaction, and cursor features.

Help

The Help menu item provides the full Konsole handbook (if KDE handbooks were installed in your Linux distribution), a “tip of the day” feature that shows interesting little-known shortcuts and tips each time you start Konsole, and the standard About Konsole dialog box.

The GNOME Terminal

As you would expect, the GNOME desktop project has its own terminal emulation program. The GNOME Terminal software package has many of the same features as Konsole and xterm. This section walks you through the various parts of configuring and using GNOME Terminal.

The Command Line Parameters

The GNOME Terminal application also provides a wealth of command line parameters that allow you to control GNOME's behavior when starting it. [Table 2.4](#) lists the parameters available.

Table 2.4 The GNOME Terminal Command Line Parameters

Parameter	Description
-e <code>command</code>	Execute the argument inside a default terminal window.
-x	Execute the entire contents of the command line after this parameter inside a default terminal window.
--window	Open a new window with a default terminal window. You may add multiple <code>--window</code> parameters to start multiple windows.
--window-with-profile= <code>profile</code>	Open a new window with a specified profile. You may also add this parameter multiple times to the command line.
--tab	Open a new tabbed terminal inside the last opened terminal window.
--tab-with-profile= <code>profile</code>	Open a new tabbed terminal inside the last opened terminal window using the specified profile.
--role= <code>role</code>	Set the role for the last specified window.
--show-menubar	Enable the menu bar at the top of the terminal window.
--hide-menubar	Disable the menu bar at the top of the terminal window.
--full-screen	Display the terminal window fully maximized.
--geometry= <code>geometry</code>	Specify the XWindow geometry parameter.

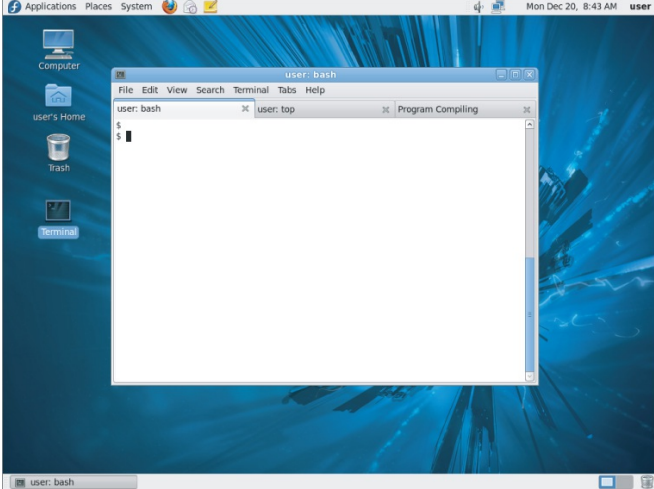
<code>--disable-factory</code>	Don't register with the activation nameserver.
<code>--use-factory</code>	Register with the activation nameserver.
<code>--startup-id=</code>	Set the ID for the Linux startup notification protocol.
<code>-t, --title=</code>	Set the window title for the terminal window.
<code>--working-directory=</code>	Set the default working directory for the terminal window.
<code>--zoom=</code>	Set the terminal's zoom factor.
<code>--active</code>	Set the last specified terminal tab as the active tab.

The GNOME Terminal command line parameters allow you to set lots of features automatically as GNOME Terminal starts. However, you can also set most of these features from within the GNOME Terminal window after it starts.

Tabs

Similar to Konsole, the GNOME Terminal calls each session a *tab*, and it also uses tabs to keep track of multiple sessions running within the window. [Figure 2.9](#) shows a GNOME Terminal window with three session tabs active.

[Figure 2.9](#) The GNOME Terminal with three active sessions



You can right-click in the tab window to see the quick menu. This quick menu provides a few actions for your use in the tab session:

- **Open Terminal:** Open a new GNOME Terminal window with a default tab session.
- **Open Tab:** Open a new session tab in the existing GNOME Terminal window.
- **Close Tab or Close Window:** If multiple tabs are open, the menu option Close Tab is shown and it closes the current session tab. If only one tab is open, the menu option Close Window is displayed and it

closes the GNOME Terminal window.

- **Copy:** Copy highlighted text in the current session tab to the clipboard.
- **Paste:** Paste data in the clipboard into the current session tab at the current cursor location.
- **Profiles:** Change the profile for the current session tab or edit the current tab profile.
- **Show Menubar:** Toggles on/off the menu bar display.
- **Input Methods:** Allows you to change the current Input Method to another character translation or turn it off completely.

The quick menu provides easy access to commonly used actions that are available from the standard menu bar in the terminal window.

The Menu Bar

The main operation of GNOME Terminal happens in the menu bar. The menu bar contains all of the configuration and customization options you'll need to make your GNOME Terminal just the way you want it. The following sections describe the different items in the menu bar.

File

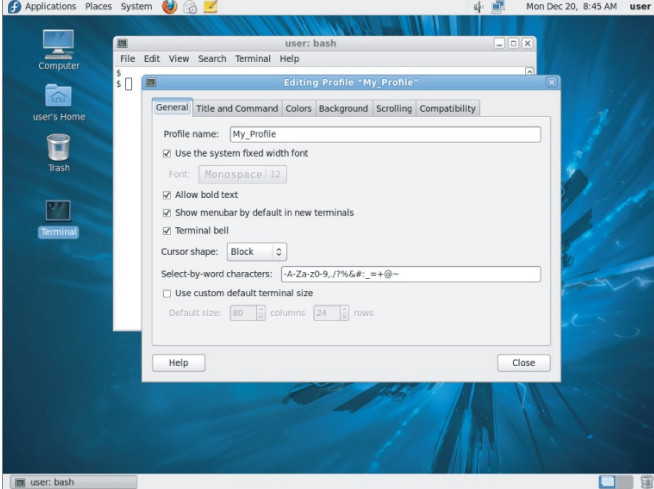
The File menu item contains items to create and manage the terminal tabs:

- **Open Terminal:** Start a new shell session in a new GNOME Terminal window.
- **Open Tab:** Start a new shell session on a new tab in the existing GNOME Terminal window.
- **New Profile:** Allows you to customize the tab session and save it as a profile, which you can recall for use later.
- **Save Contents:** Saves the contents of the scrollback buffer to a text file.
- **Close Tab:** Close the current tab in the window.
- **Close Window:** Close the current GNOME Terminal session, closing all active tabs.

Most of the items in the File menu are also available by right-clicking in the session tab area. The New Profile entry allows you to customize your session tab settings and save them for future use.

The New Profile first requests that you provide a name for the new profile; then it produces the Editing Profile dialog box, shown in [Figure 2.10](#).

[Figure 2.10](#) The *GNOME* Terminal Editing Profile dialog box



This is the area where you can set the terminal emulation features for the session. It consists of six areas:

- **General:** Provides general settings such as font, the bell, and the menu bar.
- **Title and Command:** Allows you to set the title for the session tab (displayed on the tab) and determine if the session starts with a special command rather than a shell.
- **Colors:** Sets the foreground and background colors used in the session tab.
- **Background:** Allows you to set a background image

for the session tab, or make it transparent so you can see the desktop through the session tab.

- **Scrolling:** Controls whether a scroll region is created, and how large.
- **Compatibility:** Allows you to set which control codes the Backspace and Delete keys send to the system.

Once you configure a profile, you can specify it when opening new session tabs.

Edit

The Edit menu item contains items for handling text within the tabs. You can use your mouse to copy and paste texts anywhere within the tab window. This allows you to easily copy text from the command line output to a clipboard and import it into an editor. You can also paste text from another GNOME application into the tab session.

- **Copy:** Copy selected text to the GNOME clipboard.
- **Paste:** Paste text from the GNOME clipboard into the tab session.
- **Select All:** Selects output in the entire scrollbar buffer.
- **Profiles:** Add, delete, or modify profiles in the GNOME Terminal.
- **Keyboard Shortcuts:** Create key combinations to quickly access GNOME Terminal features.
- **Profile Preferences:** Provides a quick way to edit the profile used for the current session tab.

The profile-editing feature is an extremely powerful tool for customizing several profiles, and then changing profiles as you change sessions.

View

The View menu item contains items for controlling how the session tab windows appear. They include:

- **Show Menubar:** Toggles on/off the menu bar display.
- **Full Screen:** Enlarges the GNOME Terminal window to the entire desktop.
- **Zoom In:** Makes the font in the tab window larger.
- **Zoom Out:** Makes the font in the tab window smaller.
- **Normal Size:** Returns the tab font to the default size.

If you hide the menu bar, you can easily get it back by right-clicking in any session tab and toggling the Show Menubar item.

Terminal

The Terminal menu item contains items for controlling the terminal emulation features of the tab session. They include:

- **Change Profile:** Allows you to switch to another configured profile in the session tab.
- **Set Title:** Sets the title on the session tab to easily identify it.

- **Set Character Encoding:** Selects the character set used to send and display characters.
- **Reset:** Sends the reset control code to the Linux system.
- **Reset and Clear:** Sends the reset control code to the Linux system and clears any text currently showing in the tab area.
- **Window Size List:** Lists different sizes to which the current GNOME terminal window can be adjusted. Select a size and the window automatically adjusts its size.

The character encoding offers a large list of available character sets to choose from. This is especially handy if you must work in a language other than English.

Tabs

The Tabs menu item provides items for controlling the location of the tabs and selecting which tab is active. This menu only displays when you have more than one tab session open.

- **Next Tab:** Make the next tab in the list active.
- **Previous Tab:** Make the previous tab in the list active.
- **Move Tab to the Left:** Shuffle the current tab in front of the previous tab.
- **Move Tab to the Right:** Shuffle the current tab in front of the next tab.

- **Detach Tab:** Remove the tab and start a new GNOME Terminal window using this tab session.
- **The Tab list:** Lists the currently running session tabs in the terminal window. Select a tab to quickly jump to that session.

This section allows you to manage your tabs, which can come in handy if you have several tabs open at once.

Help

The Help menu item provides a full GNOME Terminal manual so that you can research individual items and features used in the GNOME Terminal.

Summary

To start learning Linux command line commands, you need access to a command line. In a world of graphical interfaces, this can sometimes be challenging. This chapter discussed different things you should consider when trying to get to the Linux command line from within a graphical desktop environment. First, the chapter covered terminal emulation and showed what features you should know about to ensure that the Linux system can properly communicate with your terminal emulation package, and display text and graphics properly.

In particular, the chapter covered three different types of terminal emulators. The xterm terminal emulator package

was the first available for Linux. It emulates both the VT102 and Tektronix 4014 terminals. The KDE desktop project created the Konsole terminal emulation package. It provides several fancy features, such as the ability to have multiple sessions in the same window, using both console and xterm sessions, with full control of terminal emulation parameters.

Finally, the chapter discussed the GNOME desktop project's GNOME Terminal emulation package. GNOME Terminal also allows multiple terminal sessions from within a single window; plus it provides a convenient way to set many terminal features.

In the next chapter, you'll start looking at the Linux command line commands. I'll walk you through the commands necessary to navigate around the Linux filesystem, and create, delete, and manipulate files.

Chapter 3

Basic bash Shell Commands

In This Chapter

- Starting the shell
- The shell prompt
- The bash manual
- Filesystem navigation
- File and directory listing
- File handling
- Directory handling
- Viewing file contents

The default shell used in all Linux distributions is the GNU bash shell. This chapter describes the basic features available in the bash shell, and walks you through how to work with Linux files and directories using the basic commands provided by the bash shell. If you're already comfortable working with files and directories in the Linux environment, feel free to skip this chapter and continue with Chapter 4 to see more advanced commands.

Starting the Shell

The GNU bash shell is a program that provides interactive access to the Linux system. It runs as a regular program, normally started whenever a user logs in to a terminal. The shell that the system starts depends on your user ID configuration.

The `/etc/passwd` file contains a list of all the system user accounts, along with some basic configuration information about each user. Here's a sample entry from a `/etc/passwd` file:

```
rich:x:501:501:Rich Blum:/home/rich:/bin/bash
```

Each entry has seven data fields, with each field separated by a colon. The system uses the data in these fields to assign specific features for the user. These fields are:

- The username
- The user's password (or a placeholder if the password is stored in another file)
- The user's system user ID number
- The user's system group ID number
- The user's full name
- The user's default home directory
- The user's default shell program

Most of these entries will be discussed in more detail in Chapter 6. For now, just pay attention to the shell program specified.

Most Linux systems use the default bash shell program when starting a command line interface (CLI) environment for the user. The bash program also uses command line parameters to modify the type of shell you can start. [Table 3.1](#) lists the command line parameters available in bash that define what type of shell to use.

Table 3.1 The bash Command Line Parameters

Parameter	Description
<code>-c string</code>	Read commands from string and process them.
<code>-r</code>	Start a restricted shell, limiting the user to the default directory.
<code>-i</code>	Start an interactive shell, allowing input from the user.
<code>-s</code>	Read commands from the standard input.

By default, when the bash shell starts, it automatically processes commands in the `.bashrc` file in the user's home directory. Many Linux distributions use this file to also load a common file that contains commands and settings for everyone on the system. This common file is normally located in the file `/etc/bashrc`. This file often sets environment variables (see Chapter 5) used in various applications.

The Shell Prompt

Once you start a terminal emulation package or log in from the Linux console, you get access to the shell CLI *prompt*. The prompt is your gateway to the shell. This is the place

where you enter shell commands.

The default prompt symbol for the bash shell is the dollar sign (\$). This symbol indicates that the shell is waiting for you to enter text. However, you can change the format of the prompt used by your shell. The different Linux distributions use different formats for the prompt. On this Ubuntu Linux system, the bash shell prompt looks like this:

```
rich@user-desktop:~$
```

On this Fedora Linux system, it looks like this:

```
[rich@testbox~]$
```

You can configure the prompt to provide basic information about your environment. The first example shows three pieces of information in the prompt:

- The username that started the shell
- The current virtual console number
- The current directory (the tilde sign is shorthand for the home directory)

The second example provides similar information, except that it uses the hostname instead of the virtual console number. There are two environment variables that control the format of the command line prompt:

- `PS1`: Controls the format of the default command line prompt
- `PS2`: Controls the format of the second-tier command line prompt

The shell uses the default `PS1` prompt for initial data entry into the shell. If you enter a command that requires

additional information, the shell displays the second-tier prompt specified by the `PS2` environment variable.

To display the current settings for your prompts, use the `echo` command:

```
rich@ user-desktop:~$ echo $PS1
${debian_chroot:+($debian_chroot)}\u@\h:\w\$
rich@ user-desktop:~$ echo $PS2
>
rich@ user-desktop:~$
```

The format of the prompt environment variables can look pretty odd. The shell uses special characters to signify elements within the command line prompt. [Table 3.2](#) shows the special characters that you can use in the prompt string.

[Table 3.2](#) Bash Shell Prompt Characters

Character	Description
\a	Bell character
\d	Date in the format “Day Month Date”
\e	ASCII escape character
\h	Local hostname
\H	Fully qualified domain hostname
\j	Number of jobs currently managed by the shell
\l	Basename of the shell's terminal device name
\n	ASCII newline character
\r	ASCII carriage return
\s	Name of the shell
\t	Current time in 24-hour HH:MM:SS format
\T	Current time in 12-hour HH:MM:SS format
\@	Current time in 12-hour am/pm format

\u	Username of the current user
\v	Version of the bash shell
\V	Release level of the bash shell
\w	Current working directory
\W	Basename of the current working directory
!\	Bash shell history number of this command
\#	Command number of this command
\\$	A dollar sign if a normal user, or a pound sign if the root user
\nnn	Character corresponding to the octal value nnn
\\	Backslash
\[Begins a control code sequence
\]	Ends a control code sequence

Notice that all of the special prompt characters begin with a backslash (\). This is what delineates a prompt character from normal text in the prompt. In the earlier example, the prompt contained both prompt characters and a normal character (the “at” sign, and the square brackets). You can create any combination of prompt characters in your prompt. To create a new prompt, just assign a new string to the `PS1` variable:

```
[rich@testbox~]$ PS1="[t][u]\$ "
```

```
[14:40:32][rich]$
```

This new shell prompt now shows the current time, along with the username. The new `PS1` definition only lasts for the duration of the shell session. When you start a new shell, the default shell prompt definition is reloaded. In Chapter 5 you'll see how you can change the default shell prompt for

all shell sessions.

The bash Manual

Most Linux distributions include an online manual for looking up information on shell commands, as well as lots of other GNU utilities included in the distribution. It is a good idea to become familiar with the manual, as it's invaluable for working with utilities, especially when you're trying to figure out various command line parameters.

The `man` command provides access to the manual pages stored on the Linux system. Entering the `man` command followed by a specific utility name provides the manual entry for that utility. [Figure 3.1](#) shows an example of looking up the manual pages for the `date` command.

[Figure 3.1](#) Displaying the manual pages for the Linux `date` command

```
Terminal
File Edit View Search Terminal Help
DATE(1) User Commands DATE(1)

NAME
date - print or set the system date and time

SYNOPSIS
date [OPTION]... [+FORMAT]
date [-u|--utc|--universal] [MMDDhhmm[[CC]YY][.ss]]

DESCRIPTION
Display the current time in the given FORMAT, or set the system date.

-d, --date=STRING
    display time described by STRING, not 'now'

-f, --file=DATEFILE
    like --date once for each line of DATEFILE

-r, --reference=FILE
    display the last modification time of FILE

-R, --rfc-2822
    output date and time in RFC 2822 format. Example: Mon, 07 Aug 2006 12:34:56 -0600

--rfc-3339=TIMESPEC
    output date and time in RFC 3339 format. TIMESPEC='date', 'seconds', or 'ns' for
    date and time to the indicated precision. Date and time components are separated
    by a single space: 2006-08-07 12:34:56-06:00

-s, --set=STRING
Manual page date(1) line 1
```

The manual page divides information about the command into separate sections, shown in [Table 3.3](#).

Table 3.3 The Linux man Page Format

Section	Description
Name	Displays the command name and a short description
Synopsis	Shows the format of the command
Description	Describes each command option
Author	Provides information on the person who developed the command
Reporting bugs	Provides information on where to report any bugs found
Copyright	Provides information on the copyright status of the

Copyright	command code
See Also	Refers you to any similar commands available

You can step through the man pages by pressing the spacebar or using the arrow keys to scroll forward and backward through the man page text (assuming that your terminal emulation package supports the arrow key functions). When you are done with the man pages, press the q key to quit.

To see information about the bash shell, look at the man pages for it using the following command:

```
$ man bash
```

This allows you to step through all of the man pages for the bash shell. This is extremely handy when building scripts, as you don't have to refer back to books or Internet sites to look up specific formats for commands. The manual is always there for you in your session.

Filesystem Navigation

As you can see from the shell prompt, when you start a shell session, you are usually placed in your home directory. Most often, you will want to break out of your home directory and explore other areas in the Linux system. This section describes how to do that using shell commands. Before we do that, however, let's take a tour of just what the Linux filesystem looks like so we know where we're going.

The Linux Filesystem

If you're new to the Linux system, you may be confused by how it references files and directories, especially if you're used to the way that the Microsoft Windows operating system does that. Before exploring the Linux system, it helps to have an understanding of how it's laid out.

The first difference you'll notice is that Linux does not use drive letters in pathnames. In the Windows world, the physical drives installed on the PC determine the pathname of the file. Windows assigns a letter to each physical disk drive, and each drive contains its own directory structure for accessing files stored on it.

For example, in Windows you may be used to seeing the filepaths such as:

```
c:\Users\Rich\Documents\test.doc.
```

This indicates that the file test.doc is located in the directory Documents, which itself is located in the directory Rich. The Rich directory is contained under the directory Users, which is located on the hard disk partition assigned the letter C (usually the first hard drive on the PC).

The Windows filepath tells you exactly which physical disk partition contains the file named test.doc. If you wanted to save a file on a flash drive, it could be, for example, designated by the J drive. You would click the icon for the J drive, which would automatically use the filepath J:\test.doc. This path indicates that the file is located at the root of the drive assigned the letter J.

This is not the method used by Linux. Linux stores files within a single directory structure, called a *virtual directory*. The virtual directory contains filepaths from all the storage devices installed on the PC, merged into a single directory structure.

The Linux virtual directory structure contains a single base directory, called the root. Directories and files beneath the root directory are listed based on the directory path used to get to them, similar to the way Windows does it.

Tip

You'll notice that Linux uses a forward slash (/) instead of a backward slash (\) to denote directories in filepaths. The backslash character in Linux denotes an escape character and causes all sorts of problems when you use it in a filepath. This may take some getting used to if you're coming from a Windows environment.

For example, the Linux filepath `/home/rich/Documents/test.doc` indicates only that the file `test.doc` is in the directory `Documents`, under the directory `rich`, which is contained in the directory `home`. It doesn't provide any information as to which physical disk on the PC the file is stored on.

The tricky part about the Linux virtual directory is how it incorporates each storage device. The first hard drive installed in a Linux PC is called the *root drive*. The root drive contains the core of the virtual directory. Everything

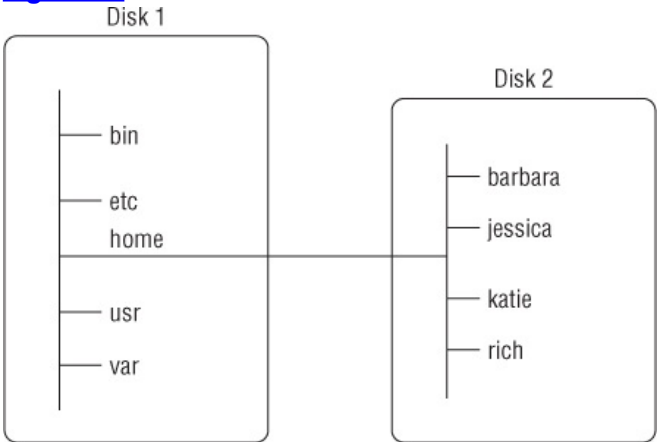
else builds from there.

On the root drive, Linux creates special directories called *mount points*. Mount points are directories in the virtual directory where you assign additional storage devices.

The virtual directory causes files and directories to appear within these mount point directories, even though they are physically stored on a different drive.

Often the system files are physically stored on the root drive, while user files are stored on a different drive, as shown in [Figure 3.2](#).

Figure 3.2 The Linux file structure



In [Figure 3.2](#), there are two hard drives on the PC. One hard drive is associated with the root of the virtual directory

(indicated by a single forward slash). Other hard drives can be mounted anywhere in the virtual directory structure. In this example, the second hard drive is mounted at the location `/home`, which is where the user directories are located.

The Linux filesystem structure has evolved from the Unix file structure. Unfortunately, the Unix file structure has been somewhat convoluted over the years by different flavors of Unix. Nowadays it seems that no two Unix or Linux systems follow the same filesystem structure. However, there are a few common directory names that are used for common functions. [Table 3.4](#) lists some of the more common Linux virtual directory names.

Table 3.4 Common Linux Directory Names

Directory	Usage
/	root of the virtual directory. Normally, no files are placed here.
/bin	binary directory, where many GNU user-level utilities are stored.
/boot	boot directory, where boot files are stored.
/dev	device directory, where Linux creates device nodes.
/etc	system configuration files directory.
/home	home directory, where Linux creates user directories.
/lib	library directory, where system and application library files are stored.
/media	media directory, a common place for mount points used for removable media.
/mnt	mount directory, another common place for mount points used for removable media.

/opt	optional directory, often used to store optional software packages.
/root	root home directory.
/sbin	system binary directory, where many GNU admin-level utilities are stored.
/tmp	temporary directory, where temporary work files can be created and destroyed.
/usr	user-installed software directory.
/var	variable directory, for files that change frequently, such as log files.

When you start a new shell prompt, your session starts in your home directory, which is a unique directory assigned to your user account. When you create a user account, the system normally assigns a unique directory for the account (see Chapter 6).

In the Windows world, you're probably used to moving around the directory structure using a graphical interface. To move around the virtual directory from a CLI prompt, you'll need to learn to use the `cd` command.

Traversing Directories

You use the change directory command (`cd`) to move your shell session to another directory in the Linux filesystem. The format of the `cd` command is pretty simplistic:

```
cd destination
```

The `cd` command may take a single parameter, `destination`, which specifies the directory name you want to

go to. If you don't specify a destination on the `cd` command, it will take you to your home directory.

The destination parameter, however, can be expressed using two different methods:

- An absolute filepath
- A relative filepath

The following sections describe the differences between these two methods.

Absolute Filepaths

You can reference a directory name within the virtual directory using an *absolute filepath*. The absolute filepath defines exactly where the directory is in the virtual directory structure, starting at the root of the virtual directory, sort of like a full name for a directory.

Thus, to reference the `apache` directory, which is contained within the `lib` directory, which in turn is contained within the `usr` directory, you would use the absolute filepath:

```
/usr/lib/NetworkManager
```

With the absolute filepath, there's no doubt as to exactly where you want to go. To move to a specific location in the filesystem using the absolute filepath, you just specify the full pathname in the `cd` command:

```
rich@testbox[~]$cd /etc
rich@testbox[etc]$
```

The prompt shows that the new directory for the shell after the `cd` command is now `/etc`. You can move to any level within the entire Linux virtual directory structure using the

absolute filepath:

```
rich@testbox[~]$ cd /usr/lib/NetworkManager  
rich@testbox[NetworkManager]$
```

However, if you're just working within your own home directory structure, often using absolute filepaths can get tedious. For example, if you're already in the directory `/home/rich`, it seems somewhat cumbersome to have to type the command

```
cd /home/rich/Documents
```

just to get to your Documents directory. Fortunately, there's a simpler solution.

Relative Filepaths

Relative filepaths allow you to specify a destination filepath relative to your current location, without having to start at the root. A relative filepath doesn't start with a forward slash, indicating the root directory.

Instead, a relative filepath starts with either a directory name (if you're traversing to a directory under your current directory), or a special character indicating a relative location to your current directory location. The two special characters used for this are:

- The dot (.) to represent the current directory
- The double dot (..) to represent the parent directory

The double dot character is extremely handy when trying to traverse a directory hierarchy. For example, if you are in the Documents directory under your home directory and

need to go to your Desktop directory, also under your home directory, you can do this:

```
rich@testbox[Documents]$ cd ../Desktop
rich@testbox[Desktop]$
```

The double dot character takes you back up one level to your home directory; then the `/Desktop` portion then takes you back down into the Desktop directory. You can use as many double dot characters as necessary to move around. For example, if you are in your home directory (`/home/rich`) and want to go to the `/etc` directory, you could type the following:

```
rich@testbox[~]$ cd ../../etc
rich@testbox[etc]$
```

Of course, in a case like this, you actually have to do more typing to use the relative filepath rather than just typing the absolute filepath, `/etc`!

File and Directory Listing

The most basic feature of the shell is the ability to see what files are available on the system. The list command (`ls`) is the tool that helps do that. This section describes the `ls` command and all of the options available to format the information it can provide.

Basic Listing

The `ls` command at its most basic form displays the files

and directories located in your current directory:

```
$ ls
4rich Desktop Download Music Pictures store store.zip test
backup Documents Drivers myprog Public store.sql Templates Videos
```

Notice that the `ls` command produces the listing in alphabetical order (in columns rather than rows). If you're using a terminal emulator that supports color, the `ls` command may also show different types of entries in different colors. The `LS_COLORS` environment variable controls this feature. Different Linux distributions set this environment variable depending on the capabilities of the terminal emulator.

If you don't have a color terminal emulator, you can use the `-F` parameter with the `ls` command to easily distinguish files from directories. Using the `-F` parameter produces the following output:

```
$ ls -F
4rich/ Documents/ Music/ Public/ store.zip Videos/
backup.zip Download/ myprog* store/ Templates/
Desktop/ Drivers/ Pictures/ store.sql test
$
```

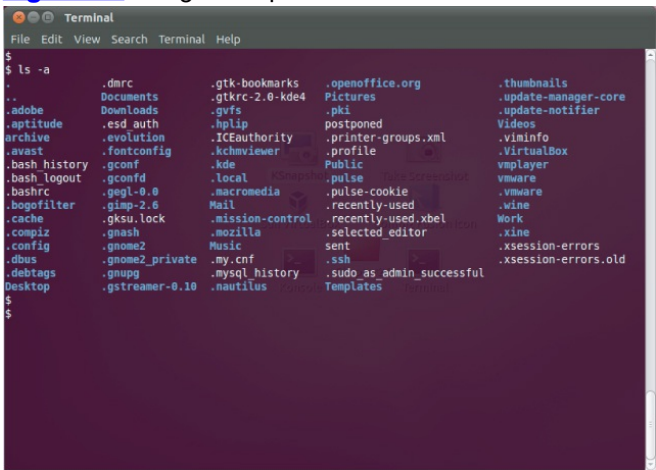
The `-F` parameter now flags the directories with a forward slash, to help identify them in the listing. Similarly, it flags executable files (like the `myprog` file above) with an asterisk, to help you more easily find the files that can be run on the system.

The basic `ls` command can be somewhat misleading. It shows the files and directories contained in the current directory, but not necessarily all of them. Linux often uses *hidden files* to store configuration information. In Linux, hidden files are files with file names that start with a period.

These files don't appear in the default `ls` listing (thus, they are called hidden).

To display hidden files along with normal files and directories, use the `-a` parameter. [Figure 3.3](#) shows an example of using the `-a` parameter with the `ls` command.

Figure 3.3 Using the `-a` parameter with the `ls` command



Wow, that's quite a difference. In a home directory for a user who has logged in to the system from a graphical desktop, you'll see lots of hidden configuration files. This particular example is from a user logged in to a GNOME desktop session. Also notice that there are three files that begin with `.bash`. These files are hidden files that are used

by the bash shell environment. These features are covered in detail in Chapter 5.

The `-R` parameter is another option the `ls` command can use. It shows files that are contained within directories in the current directory. If you have lots of directories, this can be quite a long listing. Here's a simple example of what the `-R` parameter produces:

```
$ ls -F -R
.:
file1 test1/  test2/

./test1:
myprog1*  myprog2*
./test2:
$
```

Notice that first, the `-R` parameter shows the contents of the current directory, which is a file (`file1`) and two directories (`test1` and `test2`). Following that, `-R` traverses each of the two directories, showing if any files are contained within each directory. The `test1` directory shows two files (`myprog1` and `myprog2`), while the `test2` directory doesn't contain any files. If there had been further subdirectories within the `test1` or `test2` directories, the `-R` parameter would have continued to traverse those as well. As you can see, for large directory structures this can become quite a large output listing.

Modifying the Information Presented

As you can see in the basic listings, the `ls` command doesn't produce a whole lot of information about each file.

For listing additional information, another popular parameter is `-l`. The `-l` parameter produces a long listing format, providing more information about each file in the directory:

```
$ ls -l
total 2064
drwxrwxr-x  2 rich rich    4096 2010-08-24 22:04 4rich
-rw-r--r--  1 rich rich 1766205 2010-08-24 15:34 backup.zip
drwxr-xr-x  3 rich rich    4096 2010-08-31 22:24 Desktop
drwxr-xr-x  2 rich rich    4096 2009-11-01 04:06 Documents
drwxr-xr-x  2 rich rich    4096 2009-11-01 04:06 Download
drwxrwxr-x  2 rich rich    4096 2010-07-26 18:25 Drivers
drwxr-xr-x  2 rich rich    4096 2009-11-01 04:06 Music
-rwxr--r--  1 rich rich      30 2010-08-23 21:42 myprog
drwxr-xr-x  2 rich rich    4096 2009-11-01 04:06 Pictures
drwxr-xr-x  2 rich rich    4096 2009-11-01 04:06 Public
drwxrwxr-x  5 rich rich    4096 2010-08-24 22:04 store
-rw-rw-r--  1 rich rich   98772 2010-08-24 15:30 store.sql
-rw-r--r--  1 rich rich  107507 2010-08-13 15:45 store.zip
drwxr-xr-x  2 rich rich    4096 2009-11-01 04:06 Templates
drwxr-xr-x  2 rich rich    4096 2009-11-01 04:06 Videos
[rich@testbox]$
```

The long listing format lists each file and directory contained in the directory on a single line. In addition to the file name, the listing shows additional useful information. The first line in the output shows the total number of blocks contained within the directory. Following that, each line contains the following information about each file (or directory):

- The file type—such as directory (`d`), file (`-`), character device (`c`), or block device (`b`)
- The permissions for the file (see Chapter 6)
- The number of hard links to the file (see the section “Linking Files” in this chapter)
- The username of the owner of the file

- The group name of the group the file belongs to
- The size of the file in bytes
- The time the file was modified last
- The file or directory name

The `-l` parameter is a powerful tool to have. Armed with this information, you can see just about any information you need to for any file or directory on the system.

The Complete Parameter List

There are lots of parameters for the `ls` command that can come in handy as you do file management. If you use the `man` command for `ls`, you'll see several pages of available parameters for you to use to modify the output of the `ls` command.

The `ls` command uses two types of command line parameters:

- Single-letter parameters
- Full-word (long) parameters

The single-letter parameters are always preceded by a single dash. Full-word parameters are more descriptive and are preceded by a double dash. Many parameters have both a single-letter and full-word version, while some have only one type. [Table 3.5](#) lists some of the more popular parameters that will help you out with using the bash `ls` command.

[Table 3.5](#) Some Popular ls Command Parameters

Single Letter	Full Word	Description
-a	--all	Don't ignore entries starting with a period.
-A	--almost-all	Don't list the . and .. files.
	--author	Print the author of each file.
-b	--escape	Print octal values for nonprintable characters.
	--block-size=size	Calculate the block sizes using size-byte blocks.
-B	--ignore-backups	Don't list entries with the tilde (~) symbol (used to denote backup copies).
-c		Sort by time of last modification.
-C		List entries by columns.
	--color=when	When to use colors (always, never, or auto).
-d	--directory	List directory entries instead of contents, and don't dereference symbolic links.
-F	--classify	Append file-type indicator to entries.
	--file-type	Only append file-type indicators to some filetypes (not executable files).
	--format=word	Format output as either across, commas, horizontal, long, single-column, verbose, or vertical.
-g		List full file information except for the file's owner.
	--group-directories-first	List all directories before files.
-G	--no-group	In long listing don't display group names.
-h	--human-readable	Print sizes using K for kilobytes, M for megabytes, and G for gigabytes.
	--si	Same as -h, but use powers of 1000 instead of 1024.

-i	--inode	Display the index number (inode) of each file.
-l		Display the long listing format.
-L	--dereference	Show information for the original file for a linked file.
-n	--numeric-uid-gid	Show numeric userid and groupid instead of names.
-o		In long listing don't display owner names.
-r	--reverse	Reverse the sorting order when displaying files and directories.
-R	--recursive	List subdirectory contents recursively.
-s	--size	Print the block size of each file.
-S	--sort=size	Sort the output by file size.
-t	--sort=time	Sort the output by file modification time.
-u		Display file last access time instead of last modification time.
-U	--sort=none	Don't sort the output listing.
-v	--sort=version	Sort the output by file version.
-x		List entries by line instead of columns.
-X	--sort=extension	Sort the output by file extension.

You can use more than one parameter at a time if you want to. The double dash parameters must be listed separately, but the single dash parameters can be combined together into a string behind the dash. A common combination to use is the `-a` parameter to list all files, the `-i` parameter to list the *inode* for each file, the `-l` parameter to produce a long listing, and the `-s` parameter to list the block size of the files. The inode of a file or directory

is a unique identification number the kernel assigns to each object in the filesystem. Combining all of these parameters creates the easy-to-remember `-sail` parameter:

```
$ ls -sail
total 2360
301860      8 drwx----- 36 rich rich      4096 2010-09-03 15:12 .
   65473    8 drwxr-xr-x   6 root root      4096 2010-07-29 14:20 ..
360621    8 drwxrwxr-x   2 rich rich      4096 2010-08-24 22:04 4rich
301862    8 -rw-r--r--   1 rich rich          124 2010-02-12 10:18 .bashrc
361443    8 drwxrwxr-x   4 rich rich      4096 2010-07-26 20:31 .ccache
301879    8 drwxr-xr-x   3 rich rich      4096 2010-07-26 18:25 .config
301871    8 drwxr-xr-x   3 rich rich      4096 2010-08-31 22:24 Desktop
301870    8 -rw-----   1 rich rich          26 2009-11-01 04:06 .dmrc
301872    8 drwxr-xr-x   2 rich rich      4096 2009-11-01 04:06 Download
360207    8 drwxrwxr-x   2 rich rich      4096 2010-07-26 18:25 Drivers
301882    8 drwx-----   5 rich rich      4096 2010-09-02 23:40 .gconf
301883    8 drwx-----   2 rich rich      4096 2010-09-02 23:43 .gconfd
360338    8 drwx-----   3 rich rich      4096 2010-08-06 23:06 .gftp
```

In addition to the normal `-l` parameter output information, you'll see two additional numbers added to each line. The first number in the listing is the file or directory inode number. The second number is the block size of the file. The third entry is a diagram of the type of file, along with the file's permissions. We dive into that in more detail in Chapter 6.

Following that, the next number is the number of hard links to the file (discussed later in the “Linking Files” section), the owner of the file, the group the file belongs to, the size of the file (in bytes), a timestamp showing the last modification time by default, and finally, the actual file name.

Filtering Listing Output

As you've seen in the examples, by default the `ls` command lists all of the files in a directory. Sometimes this can be overkill, especially when you're just looking for information on a single file.

Fortunately, the `ls` command also provides a way for you to define a filter on the command line. It uses the filter to determine which files or directories it should display in the output.

The filter works as a simple text-matching string. Include the filter after any command line parameters you want to use:

```
$ ls -l myprog
-rwxr--r-- 1 rich rich 30 2007-08-23 21:42 myprog
$
```

When you specify the name of specific file as the filter, the `ls` command only shows the information for that one file. Sometimes you might not know the exact name of the file you're looking for. The `ls` command also recognizes standard wildcard characters and uses them to match patterns within the filter:

- A question mark to represent one character
- An asterisk to represent zero or more characters

The question mark can be used to replace exactly one character anywhere in the filter string. For example:

```
$ ls -l mypro?
-rw-rw-r-- 1 rich rich 0 2010-09-03 16:38 myprob
-rwxr--r-- 1 rich rich 30 2010-08-23 21:42 myprog
$
```

The filter `mypro?` matched two files in the directory. Similarly, the asterisk can be used to match zero or more

characters:

```
$ ls -l myprob*  
-rw-rw-r-- 1 rich rich 0 2010-09-03 16:38 myprob  
-rw-rw-r-- 1 rich rich 0 2010-09-03 16:40 myproblem  
$
```

The asterisk matches zero characters in the `myprob` file, but it matches three characters in the `myproblem` file.

This is a powerful feature to use when searching for files when you're not quite sure of the file names.

File Handling

The bash shell provides lots of commands for manipulating files on the Linux filesystem. This section walks you through the basic commands you will need to work with files from the CLI for all your file-handling needs.

Creating Files

Every once in a while you will run into a situation where you need to create an empty file. Sometimes applications expect a log file to be present before they can write to it. In these situations, you can use the `touch` command to easily create an empty file:

```
$ touch test1  
$ ls -il test1  
1954793 -rw-r--r-- 1 rich rich 0 Sep  1 09:35 test1  
$
```

The `touch` command creates the new file you specify and assigns your username as the file owner. Because the `-il`

parameter was used for the `ls` command, the first entry in the listing shows the inode number assigned to the file. Every file on a Linux filesystem has a unique inode number.

Notice that the file size is zero because the `touch` command just created an empty file. The `touch` command can also be used to change the access and modification times on an existing file without changing the file contents:

```
$ touch test1
$ ls -l test1
-rw-r--r--  1 rich      rich      0 Sep  1 09:37 test1
$
```

The modification time of `test1` is now updated from the original time. If you want to change only the access time, use the `-a` parameter. To change only the modification time, use the `-m` parameter. By default, `touch` uses the current time. You can specify the time by using the `-t` parameter with a specific timestamp:

```
$ touch -t 201112251200 test1
$ ls -l test1
-rw-r--r--  1 rich      rich      0 Dec 25  2011 test1
$
```

Now the modification time for the file is set to a date significantly in the future from the current time.

Copying Files

Copying files and directories from one location in the filesystem to another is a common practice for system administrators. The `cp` command provides this feature.

In its most basic form, the `cp` command uses two parameters, the source object and the destination object:

cp source destination

When both the `source` and `destination` parameters are file names, the `cp` command copies the source file to a new file with the file name specified as the destination. The new file acts like a brand new file, with an updated file creation and last modified times:

```
$ cp test1 test2
$ ls -il
total 0
1954793 -rw-r--r--  1 rich    rich      0 Dec 25  2011 test1
1954794 -rw-r--r--  1 rich    rich      0 Sep  1 09:39 test2
$
```

The new file `test2` shows a different inode number, indicating that it's a completely new file. You'll also notice that the modification time for the `test2` file shows the time that it was created. If the destination file already exists, the `cp` command will prompt you to answer whether or not you want to overwrite it:

```
$ cp test1 test2
cp: overwrite 'test2'? y
$
```

If you don't answer `y`, the file copy will not proceed. You can also copy a file to an existing directory:

```
$ cp test1 dir1
$ ls -il dir1
total 0
1954887 -rw-r--r--  1 rich    rich      0 Sep  6 09:42 test1
$
```

The new file is now under the `dir1` directory, using the same file name as the original. These examples all used relative pathnames, but you can just as easily use the absolute pathname for both the source and destination objects.

To copy a file to the current directory you're in, you can use the dot symbol:

```
$ cp /home/rich/dir1/test1 .  
cp: overwrite './test1'?
```

As with most commands, the `cp` command has a few command line parameters to help you out. These are shown in [Table 3.6](#).

Table 3.6 The `cp` Command Parameters

Parameter	Description
-a	Archive files by preserving their attributes.
-b	Create a backup of each existing destination file instead of overwriting it.
-d	Preserve.
-f	Force the overwriting of existing destination files without prompting.
-i	Prompt before overwriting destination files.
-l	Create a file link instead of copying the files.
-p	Preserve file attributes if possible.
-r	Copy files recursively.
-R	Copy directories recursively.
-s	Create a symbolic link instead of copying the file.
-S	Override the backup feature.
-u	Copy the source file only if it has a newer date and time than the destination (update).
-v	Verbose mode, explaining what's happening.
-x	Restrict the copy to the current filesystem.

Use the `-p` parameter to preserve the file access or

modification times of the original file for the copied file.

```
$ cp -p test1 test3
$ ls -il
total 4
1954886 drwxr-xr-x    2 rich    rich    4096 Sep  1 09:42 dir1/
1954793 -rw-r--r--    1 rich    rich          0 Dec 25  2011 test1
1954794 -rw-r--r--    1 rich    rich          0 Sep  1 09:39 test2
1954888 -rw-r--r--    1 rich    rich          0 Dec 25  2011 test3
$
```

Now, even though the `test3` file is a completely new file, it has the same timestamps as the original `test1` file.

The `-R` parameter is extremely powerful. It allows you to recursively copy the contents of an entire directory in one command:

```
$ cp -R dir1 dir2
$ ls -l
total 8
drwxr-xr-x    2 rich    rich    4096 Sep  6 09:42 dir1/
drwxr-xr-x    2 rich    rich    4096 Sep  6 09:45 dir2/
-rw-r--r--    1 rich    rich          0 Dec 25  2011 test1
-rw-r--r--    1 rich    rich          0 Sep  6 09:39 test2
-rw-r--r--    1 rich    rich          0 Dec 25  2011 test3
$
```

Now `dir2` is a complete copy of `dir1`. You can also use wildcard characters in your `cp` commands:

```
$ cp -f test* dir2
$ ls -al dir2
total 12
drwxr-xr-x    2 rich    rich    4096 Sep  6 10:55 ./
drwxr-xr-x    4 rich    rich    4096 Sep  6 10:46 ../
-rw-r--r--    1 rich    rich          0 Dec 25  2011 test1
-rw-r--r--    1 rich    rich          0 Sep  6 10:55 test2
-rw-r--r--    1 rich    rich          0 Dec 25  2011 test3
$
```

This command copied all of the files that started with `test` to `dir2`. The `-f` parameter was included to force the overwrite of the `test1` file that was already in the directory without asking.

Linking Files

You may have noticed a couple of the parameters for the `cp` command referred to linking files. This is a pretty cool option available in the Linux filesystems. If you need to maintain two (or more) copies of the same file on the system, instead of having separate physical copies, you can use one physical copy and multiple virtual copies, called *links*. A link is a placeholder in a directory that points to the real location of the file. There are two different types of file links in Linux:

- A symbolic, or soft link
- A hard link

The hard link creates a separate file that contains information about the original file and where to locate it. When you reference the hard link file, it's just as if you're referencing the original file:

```
$ cp -l test1 test4
$ ls -il
total 16
1954886 drwxr-xr-x    2 rich    rich    4096 Sep  1 09:42 dir1/
1954889 drwxr-xr-x    2 rich    rich    4096 Sep  1 09:45 dir2/
1954793 -rw-r--r--    2 rich    rich        0 Sep  1 09:51 test1
1954794 -rw-r--r--    1 rich    rich        0 Sep  1 09:39 test2
1954888 -rw-r--r--    1 rich    rich        0 Dec 25  2011 test3
1954793 -rw-r--r--    2 rich    rich        0 Sep  1 09:51 test4
$
```

The `-l` parameter created a hard link for the `test1` file called `test4`. In the file listing, you can see that the inode number of both the `test1` and `test4` files is the same, indicating that, in reality, they are both the same file. Also

notice that the link count (the third item in the listing) now shows that both files have two links.

Note

You can only create a hard link between files on the same physical medium. You can't create a hard link between files under separate mount points. In that case, you'll have to use a soft link.

On the other hand, the `-s` parameter creates a symbolic, or soft link:

```
$ cp -s test1 test5
$ ls -il test*
total 16
1954793 -rw-r--r-- 2 rich rich 6 Sep 1 09:51 test1
1954794 -rw-r--r-- 1 rich rich 0 Sep 1 09:39 test2
1954888 -rw-r--r-- 1 rich rich 0 Dec 25 2011 test3
1954793 -rw-r--r-- 2 rich rich 6 Sep 1 09:51 test4
1954891 lrwxrwxrwx 1 rich rich 5 Sep 1 09:56 test5 -> test1
$
```

There are a couple of things to notice in the file listing. First, you'll notice that the new `test5` file has a different inode number than the `test1` file, indicating that the Linux system treats it as a separate file. Second, the file size is smaller. A linked file needs to store only information about the source file, not the actual data in the file. The file name area of the listing shows the relationship between the two files.

Tip

Instead of using the `cp` command, if you want to link files you can also use the `ln` command. By default, the `ln` command

creates hard links. If you want to create a soft link, you'll still need to use the `-s` parameter.

Be careful when copying linked files. If you use the `cp` command to copy a file that's linked to another source file, all you're doing is making another copy of the source file. This can quickly get confusing. Instead of copying the linked file, you can create another link to the original file. You can have many links to the same file with no problems. However, you also don't want to create soft links to other soft-linked files. This creates a chain of links that can not only be confusing but also be easily broken, causing all sorts of problems.

Renaming Files

In the Linux world, renaming files is called *moving*. The `mv` command is available to move both files and directories to another location:

```
$ mv test2 test6
$ ls -il test*
1954793 -rw-r--r--    2 rich      rich    6 Sep  1 09:51 test1
1954888 -rw-r--r--    1 rich      rich    0 Dec 25  2011 test3
1954793 -rw-r--r--    2 rich      rich    6 Sep  1 09:51 test4
1954891 lrwxrwxrwx    1 rich      rich    5 Sep  1 09:56 test5 -> test1
1954794 -rw-r--r--    1 rich      rich    0 Sep  1 09:39 test6
$
```

Notice that moving the file changed the file name but kept the same inode number and the timestamp value. Moving a file with soft links is a problem:

```
$ mv test1 test8
$ ls -il test*
total 16
```

```

1954888 -rw-r--r-- 1 rich rich 0 Dec 25 2011 test3
1954793 -rw-r--r-- 2 rich rich 6 Sep 1 09:51 test4
1954891 lrwxrwxrwx 1 rich rich 5 Sep 1 09:56 test5 -> test1
1954794 -rw-r--r-- 1 rich rich 0 Sep 1 09:39 test6
1954793 -rw-r--r-- 2 rich rich 6 Sep 1 09:51 test8
[rich@test2 clsc]$ mv test8 test1

```

The `test4` file that uses a hard link still uses the same inode number, which is perfectly fine. However, the `test5` file now points to an invalid file, and it is no longer a valid link.

You can also use the `mv` command to move directories:

```
$ mv dir2 dir4
```

The entire contents of the directory are unchanged. The only thing that changes is the name of the directory. Thus, the `mv` command operates much faster than the `cp` command.

Deleting Files

Most likely at some point in your Linux career, you'll want to be able to delete existing files. Whether it's to clean up a filesystem or to remove a software package, there are always opportunities to delete files.

In the Linux world, deleting is called *removing*. The command to remove files in the bash shell is `rm`. The basic form of the `rm` command is pretty simple:

```

$ rm -i test2
rm: remove 'test2'? y
$ ls -l
total 16
drwxr-xr-x 2 rich rich 4096 Sep 1 09:42 dir1/
drwxr-xr-x 2 rich rich 4096 Sep 1 09:45 dir2/
-rw-r--r-- 2 rich rich 6 Sep 1 09:51 test1
-rw-r--r-- 1 rich rich 0 Dec 25 2011 test3
-rw-r--r-- 2 rich rich 6 Sep 1 09:51 test4
lrwxrwxrwx 1 rich rich 5 Sep 1 09:56 test5 -> test1

```

\$

Notice that the command prompts you to make sure that you're serious about removing the file. There's no recycle bin or trashcan in the bash shell. Once you remove a file, it's gone forever.

Now, here's an interesting tidbit about deleting a file that has links to it:

```
$ rm test1
$ ls -l
total 12
drwxr-xr-x  2 rich  rich  4096 Sep  1 09:42 dir1/
drwxr-xr-x  2 rich  rich  4096 Sep  1 09:45 dir2/
-rw-r--r--  1 rich  rich    0 Dec 25  2011 test3
-rw-r--r--  1 rich  rich    6 Sep  1 09:51 test4
lrwxrwxrwx  1 rich  rich    5 Sep  1 09:56 test5 -> test1
$ cat test4
hello
$ cat test5
cat: test5: No such file or directory
$
```

The `test1` file was removed, which had both a hard link with the `test4` file and a soft link with the `test5` file. Notice what happened. Both of the linked files still appear, even though the `test1` file is now gone (although on my color terminal the `test5` file name now appears in red). When you look at the contents of the `test4` file that was a hard link, it still shows the contents of the file. When you look at the contents of the `test5` file that was a soft link, bash indicates that it doesn't exist anymore.

Remember that the hard link file uses the same inode number as the original file. The hard link file maintains that inode number until you remove the last file hard-linked to it, preserving the data! All the soft link file knows is that the

underlying file is now gone, so it has nothing to point to. This is an important feature to remember when working with linked files.

One other feature of the `rm` command, if you're removing lots of files and don't want to be bothered with the prompt, is to use the `-f` parameter to force the removal. Just be careful!

Tip

As with copying files, you can use wildcard characters with the `rm` command. Again, use caution when doing this, as anything you remove, even by accident, is gone forever!

Directory Handling

In Linux there are a few commands that work for both files and directories (such as the `cp` command), and some that only work for directories. To create a new directory, you'll need to use a specific command, which is covered in this section. Removing directories can get interesting, so that is covered in this section as well.

Creating Directories

There's not much to creating a new directory in Linux—just use the `mkdir` command:

```
$ mkdir dir3
```

```

ls -il
total 16
1954886 drwxr-xr-x    2 rich    rich    4096 Sep  1 09:42 dir1/
1954889 drwxr-xr-x    2 rich    rich    4096 Sep  1 10:55 dir2/
1954893 drwxr-xr-x    2 rich    rich    4096 Sep  1 11:01 dir3/
1954888 -rw-r--r--    1 rich    rich         0 Dec 25  2011 test3
1954793 -rw-r--r--    1 rich    rich         6 Sep  1 09:51 test4
$

```

The system creates a new directory and assigns it a new inode number.

Deleting Directories

Removing directories can be tricky, but there's a reason for that. There are lots of opportunities for bad things to happen when you start deleting directories. The bash shell tries to protect us from accidental catastrophes as much as possible. The basic command for removing a directory is

```

rmdir:
$ rmdir dir3
$ rmdir dir1
rmdir: dir1: Directory not empty
$

```

By default, the `rmdir` command only works for removing empty directories. Because there is a file in the `dir1` directory, the `rmdir` command refuses to remove it. You can remove nonempty directories using the `--ignore-fail-on-non-empty` parameter.

Our friend the `rm` command can also help us out some when handling directories.

If you try using it with no parameters, as with files, you'll be somewhat disappointed:

```

$ rm dir1

```

```
rm: dir1: is a directory
$
```

However, if you really want to remove a directory, you can use the `-r` parameter to recursively remove the files in the directory, then the directory itself:

```
$ rm -r dir2
rm: descend into directory 'dir2'? y
rm: remove 'dir2/test1'? y
rm: remove 'dir2/test3'? y
rm: remove 'dir2/test4'? y
rm: remove directory 'dir2'? y
$
```

While this works, it's somewhat awkward. Notice that you still must verify every file that gets removed. For a directory with lots of files and subdirectories, this can become tedious.

The ultimate solution for throwing caution to the wind and removing an entire directory, contents and all, is the `rm` command with both the `-r` and `-f` parameters:

```
$ rm -rf dir2
$
```

That's it. No warnings, no fanfare, just another shell prompt. This, of course, is an extremely dangerous tool to have, especially if you're logged in as the root user account. Use it sparingly, and only after triple checking to make sure that you're doing exactly what you want to do.

Note

You may have noticed in the last example that the two command line parameters were combined using one dash. This is a feature in the bash shell that allows you to combine command line

Viewing File Contents

So far we've covered everything there is to know about files, except for how to peek inside of them. There are several commands available for taking a look inside files without having to pull out an editor (see Chapter 11). This section demonstrates a few of the commands you have available to help you examine files.

Viewing File Statistics

You've already seen that the `ls` command can be used to provide lots of useful information about files. However, there's still more information that you can't see in the `ls` command (or at least not all at once).

The `stat` command provides a complete rundown of the status of a file on the filesystem:

```
$ stat test10
  File: "test10"
  Size: 6                Blocks: 8                Regular File
Device: 306h/774d        Inode: 1954891        Links: 2
Access: (0644/-rw-r--r--)  Uid: ( 501/ rich)  Gid: ( 501/ rich)
Access: Sat Sep  1 12:10:25 2010
Modify: Sat Sep  1 12:11:17 2010
Change: Sat Sep  1 12:16:42 2010
$
```

The results from the `stat` command show just about everything you'd want to know about the file being examined, even down to the major and minor device

numbers of the device where the file is being stored.

Viewing the File Type

Despite all of the information the `stat` command produces, there's still one piece of information missing—the file type. Before you go charging off trying to list out a 1000-byte file, it's usually a good idea to get a handle on what type of file it is. If you try listing a binary file, you'll get lots of gibberish on your monitor and possibly even lock up your terminal emulator.

The `file` command is a handy little utility to have around. It has the ability to peek inside of a file and determine just what kind of file it is:

```
$ file test1
test1: ASCII text
$ file myscript
myscript: Bourne shell script text executable
$ file myprog
myprog: ELF 32-bit LSB executable, Intel 80386, version 1 (SYSV),
dynamically linked (uses shared libs), not stripped
$
```

The `file` command classifies files into three categories:

- **Text files:** Files that contain printable characters
- **Executable files:** Files that you can run on the system
- **Data files:** Files that contain nonprintable binary characters, but that you can't run on the system

The first example shows a text file. The `file` command determined not only that the file contains text but also the character code format of the text. The second example

shows a text script file. While the file is text, because it's a script file, you can execute (run) it on the system. The final example is a binary executable program. The `file` command determines the platform that the program was compiled for and what types of libraries it requires. This is an especially handy feature if you have a binary executable program from an unknown source.

Viewing the Whole File

If you have a large text file on your hands, you may want to be able to see what's inside of it. There are three different commands in Linux that can help you out here.

The cat Command

The `cat` command is a handy tool for displaying all of the data inside a text file:

```
$ cat test1
hello

This is a test file.

That we'll use to      test the cat command.
$
```

Nothing too exciting, just the contents of the text file. There are a few parameters you can use with the `cat` command, however, that can help you out.

The `-n` parameter numbers all of the lines for you:

```
$ cat -n test1
 1  hello
 2
 3  This is a test file.
```

```
4
5
6 That we'll use to      test the cat command.
$
```

That feature will come in handy when you're examining scripts. If you just want to number the lines that have text in them, the `-b` parameter is for you:

```
$ cat -b test1
 1 hello

 2 This is a test file.

 3 That we'll use to      test the cat command.
$
```

If you need to compress multiple blank lines into a single blank line, use the `-s` parameter:

```
$ cat -s test1
hello

This is a test file.

That we'll use to      test the cat command.
$
```

Finally, if you don't want tab characters to appear, use the `-T` parameter:

```
$ cat -T test1
hello

This is a test file.

That we'll use to`Itest the cat command.
$
```

The `-T` parameter replaces any tabs in the text with the `^I` character combination.

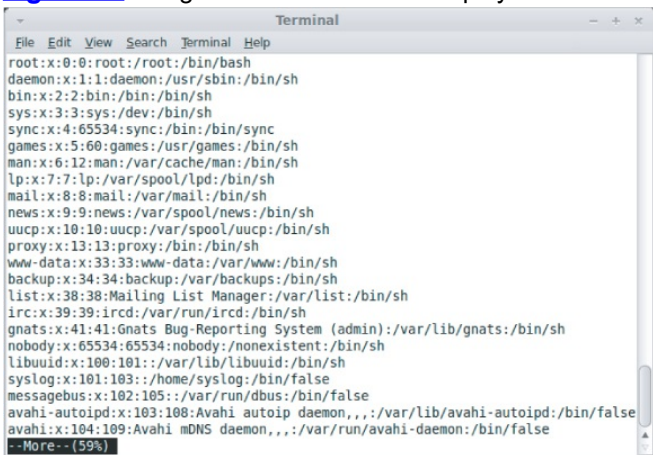
For large files, the `cat` command can be somewhat annoying. The text in the file will just quickly scroll off of the monitor without stopping. Fortunately, there's a simple way

to solve this problem.

The more Command

The main drawback of the `cat` command is that you can't control what's happening once you start it. To solve that problem, developers created the `more` command. The `more` command displays a text file, but stops after it displays each page of data. A sample `more` screen is shown in [Figure 3.4](#).

Figure 3.4 Using the `more` command to display a text file



Notice that at the bottom of the screen in [Figure 3.4](#), the `more` command displays a tag showing that you're still in the

`more` application and how far along in the text file you are. This is the prompt for the `more` command. At this point, you can enter one of several options, shown in [Table 3.7](#).

[Table 3.7](#) The `more` Command Options

Option	Description
H	Display a help menu.
spacebar	Display the next screen of text from the file.
z	Display the next screen of text from the file.
ENTER	Display one more line of text from the file.
d	Display a half-screen (11 lines) of text from the file.
q	Exit the program.
s	Skip forward one line of text.
f	Skip forward one screen of text.
b	Skip backward one screen of text.
/expression	Search for the text expression in the file.
n	Search for the next occurrence of the last specified expression.
'	Go to the first occurrence of the specified expression.
!cmd	Execute a shell command.
v	Start up the vi editor at the current line.
CTRL-L	Redraw the screen at the current location in the file.
=	Display the current line number in the file.
.	Repeat the previous command.

The `more` command allows some rudimentary movement through the text file. For more advanced features, try the `less` command.

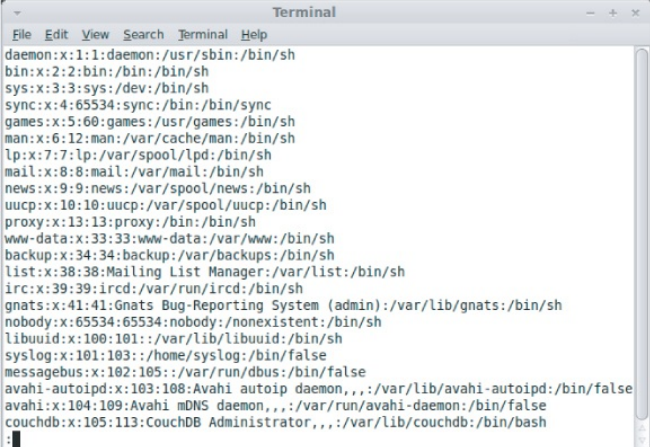
The less Command

Although from its name it sounds like it shouldn't be as advanced as the `more` command, the `less` command is actually a play on words and is an advanced version of the `more` command (the `less` command name comes from the phrase “less is more”). It provides several very handy features for scrolling both forward and backward through a text file, as well as some pretty advanced searching capabilities.

The `less` command can also display the contents of a file before it finishes reading the entire file. This is a serious drawback for both the `cat` and `more` commands when viewing extremely large files.

The `less` command operates much the same as the `more` command, displaying one screen of text from a file at a time. [Figure 3.5](#) shows the `less` command in action.

[Figure 3.5](#) Viewing a file using the `less` command



```
daemon:x:1:1:daemon:/usr/sbin:/bin/sh
bin:x:2:2:bin:/bin:/bin/sh
sys:x:3:3:sys:/dev:/bin/sh
sync:x:4:65534:sync:/bin:/bin/sync
games:x:5:60:games:/usr/games:/bin/sh
man:x:6:12:man:/var/cache/man:/bin/sh
lp:x:7:7:lp:/var/spool/lpd:/bin/sh
mail:x:8:8:mail:/var/mail:/bin/sh
news:x:9:9:news:/var/spool/news:/bin/sh
uucp:x:10:10:uucp:/var/spool/uucp:/bin/sh
proxy:x:13:13:proxy:/bin:/bin/sh
www-data:x:33:33:www-data:/var/www:/bin/sh
backup:x:34:34:backup:/var/backups:/bin/sh
list:x:38:38:Listing Manager:/var/list:/bin/sh
irc:x:39:39:ircd:/var/run/ircd:/bin/sh
gnats:x:41:41:Gnats Bug-Reporting System (admin):/var/lib/gnats:/bin/sh
nobody:x:65534:65534:nobody:/nonexistent:/bin/sh
libuuid:x:100:101::/var/lib/libuuid:/bin/sh
syslog:x:101:103::/home/syslog:/bin/false
messagebus:x:102:105::/var/run/dbus:/bin/false
avahi-autoipd:x:103:108:Avahi autoip daemon,,,:/var/lib/avahi-autoipd:/bin/false
avahi:x:104:109:Avahi mDNS daemon,,,:/var/run/avahi-daemon:/bin/false
couchdb:x:105:113:CouchDB Administrator,,,:/var/lib/couchdb:/bin/bash
```

Notice that the `less` command provides additional information in its prompt, showing the total number of lines in the file and the range of lines currently displayed. The `less` command supports the same command set as the `more` command plus lots more options. To see all of the options available, look at the man pages for the `less` command. One set of features is that the `less` command recognizes the up and down arrow keys as well as the page up and page down keys (assuming that you're using a properly defined terminal). This gives you full control when viewing a file.

Viewing Parts of a File

Often the data you want to view is located either right at the top or buried at the bottom of a text file. If the information is at the top of a large file, you still need to wait for the `cat` or `more` commands to load the entire file before you can view it. If the information is located at the bottom of a file (such as a log file), you need to wade through thousands of lines of text just to get to the last few entries. Fortunately, Linux has specialized commands to solve both of these problems.

The tail Command

The `tail` command displays the last group of lines in a file. By default, it will show the last 10 lines in the file, but you can change that with command line parameters, shown in [Table 3.8](#).

Table 3.8 The tail Command Line Parameters

Parameter	Description
<code>-c bytes</code>	Display the last byte's number of bytes in the file.
<code>-n lines</code>	Display the last line's number of lines in the file.
<code>-f</code>	Keeps the tail program active and continues to display new lines as they're added to the file.
<code>--pid=PID</code>	Along with <code>-f</code> , follows a file until the process with ID PID terminates.
<code>-s sec</code>	Along with <code>-f</code> , sleeps for sec seconds between iterations.
<code>-v</code>	Always displays output headers giving the file name.
<code>-q</code>	Never displays output headers giving the file name.

The `-f` parameter is a pretty cool feature of the `tail` command. It allows you to peek inside a file as it's being

used by other processes. The `tail` command stays active and continues to display new lines as they appear in the text file. This is a great way to monitor the system log file in real-time mode.

The head Command

While not as exotic as the `tail` command, the `head` command does what you'd expect; it displays the first group of lines at the start of a file. By default, it will display the first 10 lines of text. Similar to the `tail` command, it supports the `-c` and `-n` parameters so that you can alter what's displayed.

Usually the beginning of a file doesn't change, so the `head` command doesn't support the `-f` parameter feature. The `head` command is a handy way to just peek at the beginning of a file if you're not sure what's inside, without having to go through the hassle of displaying the entire file.

Summary

This chapter covered the basics of working with the Linux filesystem from a shell prompt. We began with a discussion of the `bash` shell and showed you how to interact with the shell. The command line interface (CLI) uses a prompt string to indicate when it's ready for you to enter commands. You can customize the prompt string to display useful information about your system, your logon ID, and even dates and times.

The bash shell provides a wealth of utilities you can use to create and manipulate files. Before you start playing with files, it's a good idea to understand how Linux stores them. This chapter discussed the basics of the Linux virtual directory and showed you how Linux references store media devices. After describing the Linux filesystem, the chapter walked you through using the `cd` command to move around the virtual directory.

After showing you how to get to a directory, the chapter demonstrated how to use the `ls` command to list the files and subdirectories. There are lots of parameters that customize the output of the `ls` command. You can obtain information on files and directories just by using the `ls` command.

The `touch` command is useful for creating empty files and for changing the access or modification times on an existing file. The chapter also discussed using the `cp` command to copy existing files from one location to another. It walked you through the process of linking files instead of copying them, providing an easy way to have the same file in two locations without making a separate copy. The `cp` command does this, as does the `ln` command.

Next, you learned how to rename files (called *moving*) in Linux using the `mv` command, and saw how to delete files (called *removing*) using the `rm` command. It also showed you how to perform the same tasks with directories, using the `mkdir` and `rmdir` commands.

Finally, the chapter closed with a discussion on viewing

the contents of files. The `cat`, `more`, and `less` commands provide easy methods for viewing the entire contents of a file, while the `tail` and `head` commands are great for peeking inside a file to just see a small portion of it.

The next chapter continues the discussion on bash shell commands. We'll take a look at more advanced administrator commands that will come in handy as you administer your Linux system.

Chapter 4

More bash Shell Commands

In This Chapter

- Managing processes
- Getting disk statistics
- Mounting new disks
- Sorting data
- Archiving data

Chapter 3 covered the basics of rummaging through the Linux filesystem and working with the files and directories. File and directory management is a major feature of the Linux shell; however, there are some more things we should look at before we start our script programming. This chapter digs into the Linux system management commands, showing you how to peek inside your Linux system using command line commands. After that, it shows you a few handy commands that you can use to work with data files on the system.

Monitoring Programs

One of the toughest jobs of being a Linux system administrator is keeping track of what's running on the system—especially now, when graphical desktops take a handful of programs just to produce a single desktop. There are always a lot of programs running on the system.

Fortunately, there are a few command line tools that can help make life easier for you. This section covers a few of the basic tools you'll need to know how to use to manage programs on your Linux system.

Peeking at the Processes

When a program runs on the system, it's referred to as a *process*. To examine these processes, you'll need to become familiar with the `ps` command, the Swiss Army knife of utilities. It can produce lots of information about all the programs running on your system.

Unfortunately, with this robustness comes complexity—in the form of numerous parameters—making the `ps` command probably one of the most difficult commands to master. Most system administrators find a subset of these parameters that provide the information they want, and then stick with using only those.

That said, however, the basic `ps` command doesn't really provide all that much information:

```
$ ps
  PID TTY          TIME CMD
 3081 pts/0    00:00:00 bash
```

Not too exciting. By default the `ps` command shows only the processes that belong to the current user and that are running on the current terminal. In this case, we only had our bash shell running (remember, the shell is just another program running on the system) and, of course, the `ps` command itself.

The basic output shows the process ID (PID) of the programs, the terminal (TTY) that they are running from, and the CPU time the process has used.

Note

The tricky feature of the `ps` command (and the part that makes it so complicated) is that at one time there were two versions of it. Each version had its own set of command line parameters controlling what information it displayed, and how. Recently, Linux developers have combined the two `ps` command formats into a single `ps` program (and of course added their own touches).

The GNU `ps` command that's used in Linux systems supports three different types of command line parameters:

- Unix-style parameters, which are preceded by a dash
- BSD-style parameters, which are not preceded by a dash
- GNU long parameters, which are preceded by a double dash

The following sections examine the three different parameter types and show examples of how they work.

Unix-Style Parameters

The Unix-style parameters originated with the original `ps` command that ran on the AT&T Unix systems invented by Bell Labs. These parameters are shown in [Table 4.1](#).

Table 4.1 The `ps` Command Unix Parameters

Parameter	Description
-A	Show all processes.
-N	Show the opposite of the specified parameters.
-a	Show all processes except session headers and processes without a terminal.
-d	Show all processes except session headers.
-e	Show all processes.
-C <i>cmslist</i>	Show processes contained in the list <i>cmslist</i> .
-G <i>grplist</i>	Show processes with a group ID listed in <i>grplist</i> .
-U <i>userlist</i>	Show processes owned by a userid listed in <i>userlist</i> .
-g <i>grplist</i>	Show processes by session or by groupid contained in <i>grplist</i> .
-p <i>pidlist</i>	Show processes with PIDs in the list <i>pidlist</i> .
-s <i>sesslist</i>	Show processes with session ID in the list <i>sesslist</i> .
-t <i>ttylist</i>	Show processes with terminal ID in the list <i>ttylist</i> .
-u <i>userlist</i>	Show processes by effective userid in the list <i>userlist</i> .
-F	Use extra full output.
-O <i>format</i>	Display specific columns in the list <i>format</i> , along with the default columns.

-M	Display security information about the process.
-c	Show additional scheduler information about the process.
-f	Display a full format listing.
-j	Show job information.
-l	Display a long listing.
-o <i>format</i>	Display only specific columns listed in <i>format</i> .
-y	Don't show process flags.
-Z	Display the security context information.
-H	Display processes in a hierarchical format (showing parent processes).
-n <i>namelist</i>	Define the values to display in the <code>WCHAN</code> column.
-w	Use wide output format, for unlimited width displays.
-L	Show process threads.
-v	Display the version of ps.

That's a lot of parameters, and remember, there are still more! The key to using the `ps` command is not to memorize all of the available parameters, only those you find most useful. Most Linux system administrators have their own sets of commonly used parameters that they remember for extracting pertinent information. For example, if you need to see everything running on the system, use the `-ef` parameter combination (the `ps` command lets you combine parameters together like this):

```
$ ps -ef
```

```

UID      PID   PPID  C  STIME TTY          TIME CMD
root         1     0  0  11:29 ?        00:00:01 init [5]
root         2     0  0  11:29 ?        00:00:00 [kthreadd]
root         3     2  0  11:29 ?        00:00:00 [migration/0]
root         4     2  0  11:29 ?        00:00:00 [ksoftirqd/0]
root         5     2  0  11:29 ?        00:00:00 [watchdog/0]
```

```

root      6      2  0 11:29 ?          00:00:00 [events/0]
root      7      2  0 11:29 ?          00:00:00 [khelper]
root     47      2  0 11:29 ?          00:00:00 [kblockd/0]
root     48      2  0 11:29 ?          00:00:00 [kacpid]
68     2349      1  0 11:30 ?          00:00:00 hald
root    2489      1  0 11:30 tty1        00:00:00 /sbin/mingetty tty1
root    2490      1  0 11:30 tty2        00:00:00 /sbin/mingetty tty2
root    2491      1  0 11:30 tty3        00:00:00 /sbin/mingetty tty3
root    2492      1  0 11:30 tty4        00:00:00 /sbin/mingetty tty4
root    2493      1  0 11:30 tty5        00:00:00 /sbin/mingetty tty5
root    2494      1  0 11:30 tty6        00:00:00 /sbin/mingetty tty6
root    2956      1  0 11:42 ?          00:00:00 /usr/sbin/httpd
apache  2958  2956  0 11:42 ?          00:00:00 /usr/sbin/httpd
apache  2959  2956  0 11:42 ?          00:00:00 /usr/sbin/httpd
root    2995      1  0 11:43 ?          00:00:00 auditd
root    2997  2995  0 11:43 ?          00:00:00 /sbin/audispd
root    3078  1981  0 12:00 ?          00:00:00 sshd: rich [priv]
rich    3080  3078  0 12:00 ?          00:00:00 sshd: rich@pts/0
rich    3081  3080  0 12:00 pts/0        00:00:00 -bash
rich    4445  3081  3 13:48 pts/0        00:00:00 ps -ef
$

```

Quite a few lines have been cut from the output to save space, but as you can see, there are lots of processes running on a Linux system. This example uses two parameters, the `-e` parameter, which shows all of the processes running on the system, and the `-f` parameter, which expands the output to show a few useful columns of information:

- **UID:** The user responsible for launching the process
- **PID:** The process ID of the process
- **PPID:** The PID of the parent process (if a process is started by another process)
- **C:** Processor utilization over the lifetime of the process
- **STIME:** The system time when the process started
- **TTY:** The terminal device from which the process

was launched

- **TIME:** The cumulative CPU time required to run the process
- **CMD:** The name of the program that was started

This produces a reasonable amount of information, which is what many system administrators would like to see. For even more information, you can use the `-l` parameter, which produces the long format output:

```
$ ps -l
 F S  UID  PID  PPID  C  PRI   NI  ADDR  SZ  WCHAN  TTY        TIME  CMD
 0 S   500 3081   3080  0   80    0  -   1173 wait pts/0    00:00:00 bash
 0 R   500 4463   3081  1   80    0  -   1116 -      pts/0    00:00:00 ps
$
```

Notice the extra columns that appear when you use the `-l` parameter:

- **F:** System flags assigned to the process by the kernel
- **S:** The state of the process (`o` = running on processor; `S` = sleeping; `R` = runnable, waiting to run; `Z` = zombie, process terminated but parent not available; `T` = process stopped)
- **PRI:** The priority of the process (higher numbers mean lower priority)
- **NI:** The nice value, which is used for determining priorities
- **ADDR:** The memory address of the process
- **SZ:** Approximate amount of swap space required if the process was swapped out
- **WCHAN:** Address of the kernel function where the process is sleeping

Before moving on, there's one more extremely handy parameter to remember, `-H`. The `-H` parameter organizes the processes in a hierarchical format, showing which processes started which other processes. Here's an extraction from an `-efH`-formatted listing:

```
$ ps -efH
UID      PID      PPID    C  STIME TTY          TIME          CMD
root    3078    1981     0  12:00 ?           00:00:00      sshd: rich [priv]
rich    3080    3078     0  12:00 ?           00:00:00      sshd: rich@pts/0
rich    3081    3080     0  12:00 pts/0       00:00:00      -bash
rich    4803    3081     1  14:31 pts/0       00:00:00      ps -efH
```

Notice the shifting in the `CMD` column output. This shows the hierarchy of the processes that are running—first, the `sshd` process started by the `root` user (this is the Secure Shell, SSH, server session, which listens for remote SSH connections). Next, because this session was connected from a remote terminal to the system, the main SSH process spawned a terminal process (`pts/0`), which in turn spawned a `bash` shell.

From there, the `ps` command was run, which appears as a child process from the `bash` process. On a multi-user system, this is a very useful tool when trying to troubleshoot runaway processes, or when trying to track down which `userid` or terminal they belong to.

BSD-Style Parameters

Now that you've seen the Unix parameters, let's take a look at the BSD-style parameters. The Berkeley Software Distribution (BSD) was a version of Unix developed at (of course) the University of California, Berkeley. It had many

subtle differences from the AT&T Unix system, thus sparking many Unix wars over the years. The BSD version of the `ps` command parameters are shown in [Table 4.2](#).

Table 4.2 The `ps` Command BSD Parameters

Parameter	Description
T	Show all processes associated with this terminal.
a	Show all processes associated with any terminal.
g	Show all processes including session headers.
r	Show only running processes.
x	Show all processes, even those without a terminal device assigned.
U <i>userlist</i>	Show processes owned by a userid listed in <i>userlist</i> .
p <i>pidlist</i>	Show processes with a PID listed in <i>pidlist</i> .
t <i>ttylist</i>	Show processes associated with a terminal listed in <i>ttylist</i> .
O <i>format</i>	List specific columns in <i>format</i> to display along with the standard columns.
x	Display data in the register format.
z	Include security information in the output.
j	Show job information.
l	Use the long format.
o <i>format</i>	Display only columns specified in <i>format</i> .
s	Use the signal format.
u	Use the user-oriented format.
v	Use the virtual memory format.
N <i>namelist</i>	Define the values to use in the <code>WCHAN</code> column.
O <i>order</i>	Define the order in which to display the information columns.

S	Sum numerical information, such as CPU and memory usage, for child processes into the parent process.
c	Display the true command name (the name of the program used to start the process).
e	Display any environment variables used by the command.
f	Display processes in a hierarchical format, showing which processes started which processes.
h	Don't display the header information.
k <i>sort</i>	Define the column(s) to use for sorting the output.
n	Use numeric values for user and group IDs, along with WCHAN information.
w	Produce wide output for wider terminals.
H	Display threads as if they were processes.
m	Display threads after their processes.
L	List all format specifiers.
v	Display the version of <code>ps</code> .

As you can see, there's a lot of overlap between the Unix and BSD types of parameters. Most of the information you can get from one you can also get from the other. Most of the time, you choose a parameter type based on which format you're more comfortable with (for example, if you were used to a BSD environment before using Linux).

When you use the BSD-style parameters, the `ps` command automatically changes the output to simulate the BSD format. Here's an example using the `l` parameter:

```
$ ps l
 F  UID   PID  PPID  PRI   NI   VSZ   RSS  WCHAN  STAT  TTY      TIME COMMAND
 0   500  3081  3080   20    0  4692  1432  wait   Ss    pts/0    0:00 -bash
 0   500  5104  3081   20    0  4468   844  -      R+    pts/0    0:00 ps l
$
```

Notice that while many of the output columns are the same as when we used the Unix-style parameters, there are a couple of different ones:

- **VSZ:** The size in kilobytes of the process in memory
- **RSS:** The physical memory that a process has used that isn't swapped out
- **STAT:** A two-character state code representing the current process state

Many system administrators like the BSD-style `l` parameter because it produces a more detailed state code for processes (the STAT column). The two-character code more precisely defines exactly what's happening with the process than the single-character Unix-style output.

The first character uses the same values as the Unix-style `s` output column, showing when a process is sleeping, running, or waiting. The second character further defines the process's status:

- `<`: The process is running at high priority.
- `N`: The process is running at low priority.
- `L`: The process has pages locked in memory.
- `s`: The process is a session leader.
- `l`: The process is multi-threaded.
- `+`: The process is running in the foreground.

From the simple example shown previously, you can see that the `bash` command is sleeping, but it is a session leader (it's the main process in my session), whereas the `ps` command was running in the foreground on the system.

The GNU Long Parameters

Finally, the GNU developers put their own touches on the new, improved `ps` command by adding a few more options to the parameter mix. Some of the GNU long parameters copy existing Unix- or BSD-style parameters, while others provide new features. [Table 4.3](#) lists the available GNU long parameters.

Table 4.3 The `ps` Command GNU Parameters

Parameter	Description
<code>--deselect</code>	Show all processes except those listed in the command line.
<code>--Group</code> <code>grplist</code>	Show processes whose group ID is listed in <code>grplist</code> .
<code>--User</code> <code>userlist</code>	Show processes whose user ID is listed in <code>userlist</code> .
<code>--group</code> <code>grplist</code>	Show processes whose effective group ID is listed in <code>grplist</code> .
<code>--pid</code> <code>pidlist</code>	Show processes whose process ID is listed in <code>pidlist</code> .
<code>--ppid</code> <code>pidlist</code>	Show processes whose parent process ID is listed in <code>pidlist</code> .
<code>--sid</code> <code>sidlist</code>	Show processes whose session ID is listed in <code>sidlist</code> .
<code>--tty</code> <code>ttylist</code>	Show processes whose terminal device ID is listed in <code>ttylist</code> .
<code>--user</code> <code>userlist</code>	Show processes whose effective user ID is listed in <code>userlist</code> .
<code>--format</code>	Display each column specified in the <code>format</code>

<code>format</code>	Display only columns specified in the <i>format</i> .
<code>--context</code>	Display additional security information.
<code>--cols n</code>	Set screen width to <i>n</i> columns.
<code>--columns n</code>	Set screen width to <i>n</i> columns.
<code>--cumulative</code>	Include stopped child process information.
<code>--forest</code>	Display processes in a hierarchical listing showing parent processes.
<code>--headers</code>	Repeat column headers on each page of output.
<code>--no-headers</code>	Don't display column headers.
<code>--lines n</code>	Set the screen height to <i>n</i> lines.
<code>--rows n</code>	Set the screen height to <i>n</i> rows.
<code>--sort order</code>	Define the column(s) to use for sorting the output.
<code>--width n</code>	Set the screen width to <i>n</i> columns.
<code>--help</code>	Display the help information.
<code>--info</code>	Display debugging information.
<code>--version</code>	Display the version of the <code>ps</code> program.

You can combine GNU long parameters with either Unix- or BSD-style parameters to really customize your display. One cool feature of GNU long parameters that we really like is the `--forest` parameter. It displays the hierarchical process information, but using ASCII characters to draw cute charts:

```

1981 ?          00:00:00 sshd
3078 ?          00:00:00  \_ sshd
3080 ?          00:00:00    \_ sshd
3081 pts/0      00:00:00      \_ bash
16676 pts/0    00:00:00        \_ ps
```

This format makes tracing child and parent processes a snap!

Real-Time Process Monitoring

The `ps` command is great for gleaning information about processes running on the system, but it has one drawback. The `ps` command can only display information for a specific point in time. If you're trying to find trends about processes that are frequently swapped in and out of memory, it's hard to do that with the `ps` command.

Instead, the `top` command can solve this problem. The `top` command displays process information similarly to the `ps` command, but it does it in real-time mode. [Figure 4.1](#) is a snapshot of the `top` command in action.

[Figure 4.1](#) The output of the `top` command while it is running


```

top - 16:04:38 up 1 min, 2 users, load average: 0.82, 0.52, 0.20
Tasks: 179 total, 1 running, 178 sleeping, 0 stopped, 0 zombie
Cpu(s): 0.5%us, 1.3%sy, 0.0%ni, 97.8%id, 0.3%wa, 0.0%hi, 0.0%si, 0.0%st
Mem: 1026084k total, 433076k used, 593008k free, 50440k buffers
Swap: 2781176k total, 0k used, 2781176k free, 191008k cached

```

PID	USER	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TIME+	COMMAND
952	root	20	0	35924	22m	7576	S	1	2.3	0:03.99	Xorg
1432	root	20	0	15056	1868	1516	S	0	0.2	0:00.25	prl_wmouse_d
1527	rich	20	0	78512	17m	13m	S	0	1.7	0:00.43	nautilus
1668	rich	20	0	64568	15m	11m	S	0	1.5	0:01.25	gnome-terminal
1	root	20	0	2804	1656	1200	S	0	0.2	0:00.61	init
2	root	20	0	0	0	0	S	0	0.0	0:00.00	kthreadd
3	root	RT	0	0	0	0	S	0	0.0	0:00.00	migration/0
4	root	20	0	0	0	0	S	0	0.0	0:00.01	ksoftirqd/0
5	root	RT	0	0	0	0	S	0	0.0	0:00.00	watchdog/0
6	root	RT	0	0	0	0	S	0	0.0	0:00.00	migration/1
7	root	20	0	0	0	0	S	0	0.0	0:00.00	ksoftirqd/1
8	root	RT	0	0	0	0	S	0	0.0	0:00.00	watchdog/1
9	root	20	0	0	0	0	S	0	0.0	0:00.00	events/0
10	root	20	0	0	0	0	S	0	0.0	0:00.04	events/1
11	root	20	0	0	0	0	S	0	0.0	0:00.00	cpuset
12	root	20	0	0	0	0	S	0	0.0	0:00.00	khelper
13	root	20	0	0	0	0	S	0	0.0	0:00.00	netns
14	root	20	0	0	0	0	S	0	0.0	0:00.00	async/mgr
15	root	20	0	0	0	0	S	0	0.0	0:00.00	pm
17	root	20	0	0	0	0	S	0	0.0	0:00.00	sync_supers

The first section of the output shows general system information. The first line shows the current time, how long the system has been up, the number of users logged in, and the load average on the system.

The load average appears as three numbers, the 1-minute, 5-minute, and 15-minute load averages. The higher the values, the more load the system is experiencing. It's not uncommon for the 1-minute load value to be high for short bursts of activity. If the 15-minute load value is high, your system may be in trouble.

Note

The trick in Linux system administration is defining what exactly a high load average value is. This value depends on what's normally running on your system and the hardware configuration. What's high for one system might be normal for another. Usually, if your load averages start getting over 2, things are getting busy on your system.

The second line shows general process information (called *tasks* in `top`): how many processes are running, sleeping, stopped, and zombie (have finished but their parent process hasn't responded).

The next line shows general CPU information. The `top` display breaks down the CPU utilization into several categories depending on the owner of the process (user versus system processes) and the state of the processes (running, idle, or waiting).

Following that, there are two lines that detail the status of the system memory. The first line shows the status of the physical memory in the system, how much total memory there is, how much is currently being used, and how much is free. The second memory line shows the status of the swap memory area in the system (if any is installed), with the same information.

Finally, the next section shows a detailed list of the currently running processes, with some information columns that should look familiar from the `ps` command output:

- **PID:** The process ID of the process

- **USER:** The user name of the owner of the process
- **PR:** The priority of the process
- **NI:** The nice value of the process
- **VIRT:** The total amount of virtual memory used by the process
- **RES:** The amount of physical memory the process is using
- **SHR:** The amount of memory the process is sharing with other processes
- **S:** The process status (D = interruptible sleep, R = running, S = sleeping, T = traced or stopped, or Z = zombie)
- **%CPU:** The share of CPU time that the process is using
- **%MEM:** The share of available physical memory the process is using
- **TIME+:** The total CPU time the process has used since starting
- **COMMAND:** The command line name of the process (program started)

By default, when you start `top` it sorts the processes based on the `%CPU` value. You can change the sort order by using one of several interactive commands while `top` is running. Each interactive command is a single character you can press while `top` is running and changes the behavior of the program. These commands are shown in [Table 4.4](#).

Table 4.4 The top Interactive Commands

Command	Description
l	Toggle the single CPU and Symmetric Multiprocessor (SMP) state.
b	Toggle the bolding of important numbers in the tables.
I	Toggle Irix/Solaris mode.
z	Configure color for the table.
l	Toggle the displaying of the load average information line.
t	Toggle the displaying of the CPU information line.
m	Toggle the displaying of the MEM and SWAP information lines.
f	Add or remove different information columns.
o	Change the display order of information columns.
F or o	Select a field on which to sort the processes (<code>%CPU</code> by default).
< or >	Move the sort field one column left (<) or right (>).
r	Toggle the normal or reverse sort order.
h	Toggle the showing of threads.
c	Toggle the showing of the command name or the full command line (including parameters) of processes.
i	Toggle the showing of idle processes.
S	Toggle the showing of the cumulative CPU time or relative CPU time.
x	Toggle highlighting of the sort field.
y	Toggle highlighting of running tasks.
z	Toggle color and mono mode.
b	Toggle bold mode for x and y modes.
u	Show processes for a specific user.

n or #	Set the number of processes to display.
k	Kill a specific process (only if process owner or if root user).
r	Change the priority (<i>renice</i>) of a specific process (only if process owner or if root user).
d or s	Change the update interval (default three seconds).
w	Write current settings to a configuration file.
q	Exit the top command.

You have lots of control over the output of the `top` command. Using this tool, you can often find offending processes that have taken over your system. Of course, once you find one, the next job is to stop it, which brings us to the next topic.

Stopping Processes

A crucial part of being a system administrator is knowing when and how to stop a process. Sometimes a process gets hung up and just needs a gentle nudge to either get going again or stop. Other times, a process runs away with the CPU and refuses to give it up. In both cases, you need a command that will allow you to control a process. Linux follows the Unix method of interprocess communication.

In Linux, processes communicate with each other using *signals*. A process signal is a predefined message that processes recognize and may choose to ignore or act on. The developers program how a process handles signals. Most well-written applications have the ability to receive and act on the standard Unix process signals. These

signals are shown in [Table 4.5](#).

Table 4.5 Linux Process Signals

Signal	Name	Description
1	HUP	Hang up.
2	INT	Interrupt.
3	QUIT	Stop running.
9	KILL	Unconditionally terminate.
11	SEGV	Segment violation.
15	TERM	Terminate if possible.
17	STOP	Stop unconditionally, but don't terminate.
18	TSTP	Stop or pause, but continue to run in background.
19	CONT	Resume execution after STOP or TSTP.

There are two commands available in Linux that allow you to send process signals to running processes.

The kill Command

The `kill` command allows you to send signals to processes based on their process ID (PID). By default, the `kill` command sends a `TERM` signal to all the PIDs listed on the command line. Unfortunately, you can only use the process PID instead of its command name, making the `kill` command difficult to use sometimes.

To send a process signal, you must either be the owner of the process or be logged in as the root user.

```
$ kill 3940
-bash: kill: (3940) - Operation not permitted
$
```

The `TERM` signal tells the process to kindly stop running. Unfortunately, if you have a runaway process, most likely it will ignore the request. When you need to get forceful, the `-s` parameter allows you to specify other signals (either using their name or signal number).

The generally accepted procedure is to first try the `TERM` signal. If the process ignores that, try the `INT` or `HUP` signals. If the program recognizes these signals, it will try to gracefully stop doing what it was doing before shutting down. The most forceful signal is the `KILL` signal. When a process receives this signal, it immediately stops running. This can lead to corrupt files.

As you can see from the following example, there's no output associated with the `kill` command.

```
# kill -s HUP 3940
#
```

To see if the command was effective, you'll have to perform another `ps` or `top` command to see if the offending process stopped.

The killall Command

The `killall` command is a powerful way to stop processes by using their names rather than the PID numbers. The `killall` command allows you to use wildcard characters as well, making it a very useful tool when you've got a system that's gone awry:

```
# killall http*
#
```

This example will kill all of the processes that start with http, such as the httpd services for the Apache Web server.

Caution

Be extremely careful using the `killall` command when logged in as the root user. It's easy to get carried away with wildcard characters and accidentally stop important system processes. This could lead to a damaged filesystem.

Monitoring Disk Space

Another important task of the system administrator is to keep track of the disk usage on the system. Whether you're running a simple Linux desktop or a large Linux server, you'll need to know how much space you have for your applications.

There are a few command line commands you can use to help you manage the media environment on your Linux system. This section describes the core commands you'll likely run into during your system administration duties.

Mounting Media

As discussed in Chapter 3, the Linux filesystem combines all media disks into a single virtual directory. Before you can use a new media disk on your system, you need to place it in the virtual directory. This task is called *mounting*.

In today's graphical desktop world, most Linux distributions have the ability to automatically mount specific types of *removable media*. A removable media device is a medium that (obviously) can be easily removed from the PC, such as CD-ROMs, floppy disks, and USB memory sticks.

If you're not using a distribution that automatically mounts and unmounts removable media, you'll have to do it yourself. This section describes the Linux command line commands to help you manage your removable media devices.

The mount Command

Oddly enough, the command used to mount media is called `mount`. By default, the `mount` command displays a list of media devices currently mounted on the system:

```
$ mount
/dev/mapper/VolGroup00-LogVol00 on / type ext3 (rw)
proc on /proc type proc (rw)
sysfs on /sys type sysfs (rw)
devpts on /dev/pts type devpts (rw,gid=5,mode=620)
/dev/sdal on /boot type ext3 (rw)
tmpfs on /dev/shm type tmpfs (rw)
none on /proc/sys/fs/binfmt_misc type binfmt_misc (rw)
sunrpc on /var/lib/nfs/rpc_pipefs type rpc_pipefs (rw)
/dev/sdbl on /media/disk type vfat
(rw,nosuid,nodev,uhelper=hal,shortname=lower,uid=503)
$
```

The `mount` command provides four pieces of information:

- The device file name of the media
- The mount point in the virtual directory where the media is mounted

- The filesystem type
- The access status of the mounted media

The last entry in the preceding example is a USB memory stick that the GNOME desktop automatically mounted at the `/media/disk` mount point. The `vfat` filesystem type shows that it was formatted on a Microsoft Windows PC.

To manually mount a media device in the virtual directory, you'll need to be logged in as the root user. The following is the basic command for manually mounting a media device:

```
mount -t type device directory
```

The `type` parameter defines the filesystem type the disk was formatted under. There are lots and lots of different filesystem types that Linux recognizes. If you share removable media devices with your Windows PCs, the types you're most likely to run into are:

- **vfat:** Windows long filesystem.
- **ntfs:** Windows advanced filesystem used in Windows NT, XP, and Vista.
- **iso9660:** The standard CD-ROM filesystem.

Most USB memory sticks and floppies are formatted using the `vfat` filesystem. If you need to mount a data CD, you'll have to use the `iso9660` filesystem type.

The next two parameters define the location of the device file for the media device and the location in the virtual directory for the mount point. For example, to manually mount the USB memory stick at device `/dev/sdb1` at location

/media/disk, you'd use the following command:

```
mount -t vfat /dev/sdb1 /media/disk
```

Once a media device is mounted in the virtual directory, the root user will have full access to the device, but access by other users will be restricted. You can control who has access to the device using directory permissions (discussed in Chapter 6).

In case you need to use some of the more exotic features of the `mount` command, the available parameters are shown in [Table 4.6](#).

Table 4.6 The mount Command Parameters

Parameter	Description
-a	Mount all filesystems specified in the <code>/etc/fstab</code> file.
-f	Causes the mount command to simulate mounting a device, but not actually mount it.
-F	When used with the <code>-a</code> parameter, mounts all filesystems at the same time.
-v	Verbose mode, explains all the steps required to mount the device.
-I	Don't use any filesystem helper files under <code>/sbin/mount.filesystem.</code>
-l	Add the filesystem labels automatically for ext2, ext3, or XFS filesystems.
-n	Mount the device without registering it in the <code>/etc/mstab</code> mounted device file.
-p <i>num</i>	For encrypted mounting, read the passphrase from the file descriptor <i>num</i> .
-s	Ignore mount options not supported by the filesystem.

-r	Mount the device as read-only.
-w	Mount the device as read-write (the default).
-L <i>label</i>	Mount the device with the specified <i>label</i> .
-U <i>uuid</i>	Mount the device with the specified <i>uuid</i> .
-O	When used with the -a parameter, limits the set of filesystems applied.
-o	Add specific options to the filesystem.

The `-o` option allows you to mount the filesystem with a comma-separated list of additional options. The popular options to use are:

- `ro`: Mount as read-only.
- `rw`: Mount as read-write.
- `user`: Allow an ordinary user to mount the filesystem.
- `check=none`: Mount the filesystem without performing an integrity check.
- `loop`: Mount a file.

A popular thing in Linux these days is to distribute a CD as a `.iso` file. The `.iso` file is a complete image of the CD in a single file. Most CD-burning software packages can create a new CD based on the `.iso` file. A feature of the `mount` command is that you can mount a `.iso` file directly to your Linux virtual directory without having to burn it onto a CD. This is accomplished using the `-o` parameter with the `loop` option:

```
$ mkdir mnt
$ su
Password:
# mount -t iso9660 -o loop MEPIS-KDE4-LIVE-DVD_32.iso mnt
# ls -l mnt
total 16
```

```

-r--r--r-- 1 root root    702 2007-08-03 08:49 about
dr-xr-xr-x 3 root root  2048 2007-07-29 14:30 boot
-r--r--r-- 1 root root  2048 2007-08-09 22:36 boot.catalog
-r--r--r-- 1 root root    894 2004-01-23 13:22 cdrom.iso
-r--r--r-- 1 root root  5229 2006-07-07 18:07 MCWL
dr-xr-xr-x 2 root root  2048 2007-08-09 22:32 mepis
dr-xr-xr-x 2 root root  2048 2007-04-03 16:44 OSX
-r--r--r-- 1 root root    107 2007-08-09 22:36 version

# cd mnt/boot
# ls -l
total 4399
dr-xr-xr-x 2 root root      2048 2007-06-29 09:00 grub
-r--r--r-- 1 root root 2392512 2007-07-29 12:53 initrd.gz
-r--r--r-- 1 root root   94760 2007-06-14 14:56 memtest
-r--r--r-- 1 root root 2014704 2007-07-29 14:26 vmlinuz
#

```

The `mount` command mounted the `.iso` CD image file just as if it were a real CD and allowed us to maneuver around within its filesystem.

The *umount* Command

To remove a removable media device, you should never just remove it from the system. Instead, you should always *umount* it first.

Tip

Linux doesn't allow you to eject a mounted CD. If you ever have trouble removing a CD from the drive, most likely it means the CD is still mounted in the virtual directory. Unmount it first, and then try to eject it.

The command used to unmount devices is `umount` (yes, there's no “n” in the command, which gets confusing sometimes). The format for the `umount` command is pretty

simple:

```
umount [directory | device ]
```

The `umount` command gives you the choice of defining the media device by either its device location or its mounted directory name. If any program has a file open on a device, the system won't let you unmount it.

```
[root@testbox mnt]# umount /home/rich/mnt
umount: /home/rich/mnt: device is busy
umount: /home/rich/mnt: device is busy
[root@testbox mnt]# cd /home/rich
[root@testbox rich]# umount /home/rich/mnt
[root@testbox rich]# ls -l mnt
total 0
[root@testbox rich]#
```

In this example, the command prompt was still in a directory within the filesystem structure, so the `umount` command couldn't unmount the image file. Once the command prompt was moved out of the image file filesystem, the `umount` command was able to successfully unmount the image file.

Using the df Command

Sometimes you need to see how much disk space is available on an individual device. The `df` command allows you to easily see what's happening on all of the mounted disks:

```
$ df
```

Filesystem	1K-blocks	Used	Available	Use%	Mounted on
/dev/sda2	18251068	7703964	9605024	45%	/
/dev/sda1	101086	18680	77187	20%	/boot
tmpfs	119536	0	119536	0%	/dev/shm
/dev/sdb1	127462	113892	13570	90%	/media/disk

```
$
```

The `df` command shows each mounted filesystem that contains data. As you can see from the `mount` command earlier, some mounted devices are used for internal system purposes. The command displays the following:

- The device location of the device
- How many 1024-byte blocks of data it can hold
- How many 1024-byte blocks are used
- How many 1024-byte blocks are available
- The amount of used space as a percentage
- The mount point where the device is mounted

A few different command line parameters are available with the `df` command, most of which you'll never use. One popular parameter is `-h`, which shows the disk space in human-readable form, usually as an M for megabytes or a G for gigabytes:

```
$ df -h
Filesystem      Size  Used Avail Use% Mounted on
/dev/sdb2        18G   7.4G   9.2G   45% /
/dev/sda1        99M    19M    76M   20% /boot
tmpfs            117M     0   117M    0% /dev/shm
/dev/sdb1       125M   112M    14M   90% /media/disk
$
```

Now instead of having to decode those ugly block numbers, all of the disk sizes are shown using “normal” sizes. The `df` command is invaluable in troubleshooting disk space problems on the system.

Note

Remember that the Linux system always has processes running in the background that handle

files. The values from the `df` command reflect what the Linux system thinks are the current values at that point in time. It's possible that you have a process running that has created or deleted a file but has not released the file yet. This value is not included in the free space calculation.

Using the `du` Command

With the `df` command, it is easy to see when a disk is running out of space. The next problem for the system administrator is to know what to do when that happens.

Another useful command to help you out is the `du` command. The `du` command shows the disk usage for a specific directory (by default, the current directory). This is a quick way to determine if you have any obvious disk hogs on the system.

By default, the `du` command displays all of the files, directories, and subdirectories under the current directory, and it shows how many disk blocks each file or directory takes. For a standard-sized directory, this can be quite a listing. Here's a partial listing of using the `du` command:

```
$ du
484  ./gstreamer-0.10
8    ./Templates
8    ./Download
8    ./ccache/7/0
24   ./ccache/7
368  ./ccache/a/d
384  ./ccache/a
424  ./ccache
8    ./Public
8    ./gphpedit/plugins
32   ./gphpedit
72   ./gconfd
```



```
128      ./nautilusmetafiles
384      ./nautilus
72      ./bittorrent/data/metainfo
20      ./bittorrent/data/resume
144     ./bittorrent/data
152     ./bittorrent
8       ./Videos
8       ./Music
16      ./config/gtk-2.0
40      ./config
8       ./Documents
```

The number at the right of each line is the number of disk blocks that each file or directory takes. Notice that the listing starts at the bottom of a directory and works its way up through the files and subdirectories contained within the directory.

The `du` command by itself can be somewhat useless. It's nice to be able to see how much disk space each individual file and directory takes up, but it can be meaningless when you have to wade through pages and pages of information before you find what you're looking for.

There are a few command line parameters that you can use with the `du` command to make things a little more legible:

- `-c`: Produce a grand total of all the files listed.
- `-h`: Print sizes in human-readable form, using K for kilobyte, M for megabyte, and G for gigabyte.
- `-s`: Summarize each argument.

The next step for the system administrator is to use some file-handling commands for manipulating large amounts of data. That's exactly what the next section covers.

Working with Data Files

When you have a large amount of data, it's often difficult to handle the information and make it useful. As you saw with the `du` command in the previous section, it's easy to get data overload when working with system commands.

The Linux system provides several command line tools to help you manage large amounts of data. This section covers the basic commands that every system administrator—as well as any everyday Linux user—should know how to use to make their lives easier.

Sorting Data

One popular function that comes in handy when working with large amounts of data is the `sort` command. The `sort` command does what it says—it sorts data.

By default, the `sort` command sorts the data lines in a text file using standard sorting rules for the language you specify as the default for the session.

```
$ cat file1
one
two
three
four
five
$ sort file1
five
four
one
three
two
$
```

Pretty simple. However, things aren't always as easy as they appear. Take a look at this example:

```
$ cat file2
1
2
100
45
3
10
145
75
$ sort file2
1
10
100
145
2
3
45
75
$
```

If you were expecting the numbers to sort in numerical order, you were disappointed. By default, the `sort` command interprets numbers as characters and performs a standard character sort, producing output that might not be what you want. To solve this problem, use the `-n` parameter, which tells the `sort` command to recognize numbers as numbers instead of characters and to sort them based on their numerical values:

```
$ sort -n file2
1
2
3
10
45
75
100
145
$
```

Now, that's much better! Another common parameter that's used is `-M`, the month sort. Linux log files usually contain a timestamp at the beginning of the line to indicate when the event occurred:

```
Sep 13 07:10:09 testbox smartd[2718]: Device: /dev/sda, opened
```

If you sort a file that uses timestamp dates using the default sort, you'll get something like this:

```
$ sort file3
Apr
Aug
Dec
Feb
Jan
Jul
Jun
Mar
May
Nov
Oct
Sep
$
```

Not exactly what you wanted. If you use the `-M` parameter, the sort command recognizes the three-character month nomenclature, and sorts appropriately:

```
$ sort -M file3
Jan
Feb
Mar
Apr
May
Jun
Jul
Aug
Sep
Oct
Nov
Dec
$
```

There are several other handy `sort` parameters to use, as

shown in [Table 4.7](#).

Table 4.7 The sort Command Parameters

Single Dash	Double Dash	Description
-b	--ignore-leading-blanks	Ignore leading blanks when sorting.
-c	--check = quiet	Don't sort, but don't report if data is out of sort order.
-c	--check	Don't sort, but check if the input data is already sorted. Report if not sorted.
-d	--dictionary-order	Consider only blanks and alphanumeric characters; don't consider special characters.
-f	--ignore-case	By default, sort orders capitalized letters first. This parameter ignores case.
-g	--general-numeric-sort	Use general numerical value to sort.
-i	--ignore-nonprinting	Ignore nonprintable characters in the sort.
-k	--key = POS1[,POS2]	Sort based on position POS1, and end at POS2 if specified.
-M	--month-sort	Sort by month order using three-character month names.
-m	--merge	Merge two already sorted data files.
-n	--numeric-sort	Sort by string numerical value.
-o	-output = file	Write results to file specified.
	--random	

-R	--random-sort	Sort by a random hash of keys.
	--random-source = <i>FILE</i>	Specify the file for random bytes used by the -R parameter.
-r	--reverse	Reverse the sort order (descending instead of ascending).
-S	--buffer-size = <i>SIZE</i>	Specify the amount of memory to use.
-s	--stable	Disable last-resort comparison.
-T	--temporary-direction = <i>DIR</i>	Specify a location to store temporary working files.
-t	--field-separator = <i>SEP</i>	Specify the character used to distinguish key positions.
-u	--unique	With the -c parameter, check for strict ordering; without the -c parameter, output only the first occurrence of two similar lines.
-z	--zero-terminated	End all lines with a NULL character instead of a new line.

The `-k` and `-t` parameters are handy when sorting data that uses fields, such as the `/etc/passwd` file. Use the `-t` parameter to specify the field separator character, and the `-k` parameter to specify which field to sort on. For example, to sort the password file based on numerical userid, just do this:

```
$ sort -t ':' -k 3 -n /etc/passwd
root:x:0:0:root:/root:/bin/bash
bin:x:1:1:bin:/bin:/sbin/nologin
daemon:x:2:2:daemon:/sbin:/sbin/nologin
adm:x:3:4:adm:/var/adm:/sbin/nologin
lp:x:4:7:lp:/var/spool/lpd:/sbin/nologin
```

```
sync:x:5:0:sync:/sbin:/bin/sync
shutdown:x:6:0:shutdown:/sbin:/sbin/shutdown
halt:x:7:0:halt:/sbin:/sbin/halt
mail:x:8:12:mail:/var/spool/mail:/sbin/nologin
news:x:9:13:news:/etc/news:
uucp:x:10:14:uucp:/var/spool/uucp:/sbin/nologin
operator:x:11:0:operator:/root:/sbin/nologin
games:x:12:100:games:/usr/games:/sbin/nologin
gopher:x:13:30:gopher:/var/gopher:/sbin/nologin
ftp:x:14:50:FTP User:/var/ftp:/sbin/nologin
```

Now the data is perfectly sorted based on the third field, which is the numerical userid value.

The `-n` parameter is great for sorting numerical outputs, such as the output of the `du` command:

```
$ du -sh * | sort -nr
1008k  mrtg-2.9.29.tar.gz
972k   bldgl
888k   fbs2.pdf
760k   Printtest
680k   rsync-2.6.6.tar.gz
660k   code
516k   figl001.tiff
496k   test
496k   php-common-4.0.4pl1-6mdk.i586.rpm
448k   MesaGLUT-6.5.1.tar.gz
400k   plp
```

Notice that the `-r` option also sorts the values in descending order, so you can easily see what files are taking up the most space in your directory.

Note

The pipe command (`|`) used in this example redirects the output of the `du` command to the `sort` command. That's discussed in more detail in Chapter 10.

Searching for Data

Often in a large file, you have to look for a specific line of data buried somewhere in the middle of the file. Instead of manually scrolling through the entire file, you can let the `grep` command search for you. The command line format for the `grep` command is:

```
grep [options] pattern [file]
```

The `grep` command searches either the input or the file you specify for lines that contain characters that match the specified pattern. The output from `grep` is the lines that contain the matching pattern.

Here are two simple examples of using the `grep` command with the `file1` file used in the “Sorting Data” section:

```
$ grep three file1
three
$ grep t file1
two
three
$
```

The first example searches the file `file1` for text matching the pattern *three*. The `grep` command produces the line that contains the matching pattern. The next example searches the file `file1` for the text matching the pattern *t*. In this case, there were two lines that matched the specified pattern, and both are displayed.

Because of the popularity of the `grep` command, it has undergone lots of development changes over its lifetime. There are lots of features that have been added to the `grep` command. If you look over the man pages for the `grep`

command, you'll see how versatile it really is.

If you want to reverse the search (output lines that don't match the pattern) use the `-v` parameter:

```
$ grep -v t file1
one
four
five
$
```

If you need to find the line numbers where the matching patterns are found, use the `-n` parameter:

```
$ grep -n t file1
2:two
3:three
$
```

If you just need to see a count of how many lines contain the matching pattern, use the `-c` parameter:

```
$ grep -c t file1
2
$
```

If you need to specify more than one matching pattern, use the `-e` parameter to specify each individual pattern:

```
$ grep -e t -e f file1
two
three
four
five
$
```

This example outputs lines that contain either the string `t` or the string `f`.

By default, the `grep` command uses basic Unix-style regular expressions to match patterns. A Unix-style regular expression uses special characters to define how to look for matching patterns.

For a more detailed explanation of regular expressions, see Chapter 19.

Here's a simple example of using a regular expression in a `grep` search:

```
$ grep [tf] file1
two
three
four
five
$
```

The square brackets in the regular expression indicate that `grep` should look for matches that contain either a *t* or an *f* character. Without the regular expression, `grep` would search for text that would match the string `tf`.

The `egrep` command is an offshoot of `grep`, which allows you to specify POSIX extended regular expressions, which contain more characters for specifying the matching pattern (again, see Chapter 19 for more details). The `fgrep` command is another version that allows you to specify matching patterns as a list of fixed-string values, separated by newline characters. This allows you to place a list of strings in a file and then use that list in the `fgrep` command to search for the strings in a larger file.

Compressing Data

If you've done any work in the Microsoft Windows world, no doubt you've used zip files. It became such a popular feature that Microsoft eventually incorporated it into the Windows XP operating system. The zip utility allows you to easily compress large files (both text and executable) into

smaller files that take up less space.

Linux contains several file compression utilities. While this may sound great, it often leads to confusion and chaos when trying to download files. [Table 4.8](#) lists the file compression utilities available for Linux.

Table 4.8 Linux File Compression Utilities

Utility	File Extension	Description
<code>bzip2</code>	<code>.bz2</code>	Uses the Burrows-Wheeler block sorting text compression algorithm and Huffman coding
<code>compress</code>	<code>.Z</code>	Original Unix file compression utility, starting to fade away into obscurity
<code>gzip</code>	<code>.gz</code>	The GNU Project's compression utility, uses Lempel-Ziv coding
<code>zip</code>	<code>.zip</code>	The Unix version of the PKZIP program for Windows

The `compress` file compression utility is not often found on Linux systems. If you download a file with a `.Z` extension, you can usually install the `compress` package (called `ncompress` in many Linux distributions) using the software installation methods discussed in Chapter 8, and then uncompress the file with the `uncompress` command.

The *bzip2* Utility

The `bzip2` utility is a relatively new compression package that is gaining popularity, especially when compressing large binary files. The utilities in the `bzip2` package are:

- `bzip2` for compressing files

- `bzcat` for displaying the contents of compressed text files
- `bunzip2` for uncompressing compressed `.bz2` files
- `bzip2recover` for attempting to recover damaged compressed files

By default, the `bzip2` command attempts to compress the original file and replaces it with the compressed file, using the same file name with a `.bz2` extension:

```
$ ls -l myprog
-rwxrwxr-x 1 rich rich 4882 2007-09-13 11:29 myprog
$ bzip2 myprog
$ ls -l my*
-rwxrwxr-x 1 rich rich 2378 2007-09-13 11:29 myprog.bz2
$
```

The original size of the `myprog` program was 4882 bytes, and after the `bzip2` compression it is now 2378 bytes. Also, notice that the `bzip2` command automatically renamed the original file with the `.bz2` extension, indicating what compression technique we used to compress it.

To uncompress the file, just use the `bunzip2` command:

```
$ bunzip2 myprog.bz2
$ ls -l myprog
-rwxrwxr-x 1 rich rich 4882 2007-09-13 11:29 myprog
$
```

As you can see, the uncompressed file is back to the original file size. Once you compress a text file, you can't use the standard `cat`, `more`, or `less` commands to view the data. Instead, you need to use the `bzcat` command:

```
$ bzcat test.bz2
This is a test text file.
The quick brown fox jumps over the lazy dog.
This is the end of the test text file.
$
```

The `bzcat` command displays the text inside the compressed file without uncompressing the actual file.

The *gzip* Utility

By far the most popular file compression utility in Linux is the `gzip` utility. The `gzip` package is a creation of the GNU Project, in their attempt to create a free version of the original Unix `compress` utility. This package includes the files:

- `gzip` for compressing files
- `gzcat` for displaying the contents of compressed text files
- `gunzip` for uncompressing files

These utilities work the same way as the `bzip2` utilities:

```
$ gzip myprog
$ ls -l my*
-rwxrwxr-x 1 rich rich 2197 2007-09-13 11:29 myprog.gz
$
```

The `gzip` command compresses the file you specify on the command line. You can also specify more than one file name or even use wildcard characters to compress multiple files at once:

```
$ gzip my*
$ ls -l my*
-rwxr--r-- 1 rich rich 103 Sep 6 13:43 myprog.c.gz
-rwxr-xr-x 1 rich rich 5178 Sep 6 13:43 myprog.gz
-rwxr--r-- 1 rich rich 59 Sep 6 13:46 myscript.gz
-rwxr--r-- 1 rich rich 60 Sep 6 13:44 myscript. gz
$
```

The `gzip` command compresses every file in the directory that matches the wildcard pattern.

The zip Utility

The `zip` utility is compatible with the popular PKZIP package created by Phil Katz for MS-DOS and Windows. There are four utilities in the Linux `zip` package:

- `zip` creates a compressed file containing listed files and directories.
- `zipcloak` creates an encrypted compress file containing listed files and directories.
- `zipnote` extracts the comments from a zip file.
- `zipsplit` splits a zip file into smaller files of a set size (used for copying large zip files to floppy disks).
- `unzip` extracts files and directories from a compressed zip file.

To see all of the options available for the `zip` utility, just enter it by itself on the command line:

```
$ zip
Copyright (C) 1990-2005 Info-ZIP
Type 'zip -L' for software license.
Zip 2.31 (March 8th 2005). Usage:
zip [-options] [-b path] [-t mmdyyyyy] [-n suffixes] [zipfile list]
[-xi list]

The default action is to add or replace zipfile entries from list,
which can include the special name - to compress standard input.
If zipfile and list are omitted, zip compresses stdin to stdout.
-f freshen: only changed files      -u update: only changed or new files
-d delete entries in zipfile        -m move into zipfile (delete files)
-r recurse into directories         -j junk directory names
-o store only                       -l convert LF to CR LF
-l compress faster                  -9 compress better
-q quiet operation                  -v verbose operation
-c add one-line comments           -z add zipfile comment
-@ read names from stdin            -o make file as old as latest entry
-x exclude the following names     -i include only the following names
-F fix zipfile (-FF try harder)    -D do not add directory entries
-A adjust self-extracting exe      -J junk zipfile prefix (unzipsfx)
```

```
-T test zipfile integrity          -X eXclude eXtra file attributes
-y store symbolic links as the link instead of the referenced file
-R PKZIP recursion (see manual)
-e encrypt                        -n don't compress these suffixes
$
```

The power of the `zip` utility is its ability to compress entire directories of files into a single compressed file. This makes it ideal for archiving entire directory structures:

```
$ zip -r testzip test
adding: test/ (stored 0%)
adding: test/test1/ (stored 0%)
adding: test/test1/myprog2 (stored 0%)
adding: test/test1/myprog1 (stored 0%)
adding: test/myprog.c (deflated 39%)
adding: test/file3 (deflated 2%)
adding: test/file4 (stored 0%)
adding: test/test2/ (stored 0%)
adding: test/file1.gz (stored 0%)
adding: test/file2 (deflated 4%)
adding: test/myprog.gz (stored 0%)
$
```

This example creates the zip file named `testzip.zip` and recurses through the directory `test`, adding each file and directory found to the zip file. Notice from the output that not all of the files stored in the zip file could be compressed. The `zip` utility automatically determines the best compression type to use for each individual file.

Caution

When you use the recursion feature in the `zip` command, files are stored in the same directory structure in the zip file. Files contained in subdirectories are stored in the zip file within the same subdirectories. You must be careful when extracting the files; the `unzip` command will rebuild the entire directory structure in the new location. Sometimes this gets confusing when you have lots of subdirectories and files.

Archiving Data

While the `zip` command works great for compressing and archiving data into a single file, it's not the standard utility used in the Unix and Linux worlds. By far the most popular archiving tool used in Unix and Linux is the `tar` command.

The `tar` command was originally used to write files to a tape device for archiving. However, it can also write the output to a file, which has become a popular way to archive data in Linux.

The following is the format of the `tar` command:

```
tar function [options] object1 object2 ...
```

The function parameter defines what the `tar` command should do, as shown in [Table 4.9](#).

Table 4.9 The tar Command Functions

Function	Long Name	Description
-A	-- concatenate	Append an existing tar archive file to another existing tar archive file.
-c	--create	Create a new tar archive file.
-d	--diff	Check the differences between a tar archive file and the filesystem.
	--delete	Delete from an existing tar archive file.
-r	--append	Append files to the end of an existing tar archive file.
-t	--list	List the contents of an existing tar archive file.
		Append files to an existing tar archive file that are

-u	--update	newer than a file with the same name in the existing archive.
-x	--extract	Extract files from an existing archive file.

Each function uses *options* to define a specific behavior for the tar archive file. [Table 4.10](#) lists the common options that you can use with the `tar` command.

Table 4.10 The tar Command Options

Option	Description
-C <i>dir</i>	Change to the specified directory.
-f <i>file</i>	Output results to file (or device) <i>file</i> .
-j	Redirect output to the <code>bzip2</code> command for compression.
-p	Preserve all file permissions.
-v	List files as they are processed.
-z	Redirect the output to the <code>gzip</code> command for compression.

These options are usually combined to create the following scenarios. First, you'll want to create an archive file using this command:

```
tar -cvf test.tar test/ test2/
```

The above command creates an archive file called `test.tar` containing the contents of both the `test` directory and the `test2` directory. Next, this command:

```
tar -tf test.tar
```

lists (but doesn't extract) the contents of the tar file `test.tar`. Finally, this command:

```
tar -xvf test.tar
```

extracts the contents of the tar file `test.tar`. If the tar file was

created from a directory structure, the entire directory structure is re-created starting at the current directory.

As you can see, using the `tar` command is a simple way to create archive files of entire directory structures. This is a common method for distributing source code files for open source applications in the Linux world.

Tip

If you download open source software, often you'll see filenames that end in `.tgz`. These are gzipped tar files, and can be extracted using the command `tar -zxvf filename.tgz`.

Summary

This chapter discussed some of the more advanced bash commands used by Linux system administrators and programmers. The `ps` and `top` commands are vital in determining the status of the system, allowing you to see what applications are running and how many resources they are consuming.

In this day of removable media, another popular topic for system administrators is mounting storage devices. The `mount` command allows you to mount a physical storage device into the Linux virtual directory structure. To remove the device, use the `umount` command.

Finally, the chapter discussed various utilities used for

handling data. The `sort` utility easily sorts large data files to help you organize data, and the `grep` utility allows you to quickly scan through large data files looking for specific information. There are a few different file compression utilities available in Linux, including `bzip2`, `gzip`, and `zip`. Each one allows you to compress large files to help save space on your filesystem. The Linux `tar` utility is a popular way to archive directory structures into a single file that can easily be ported to another system.

The next chapter discusses Linux environment variables. Environment variables allow you to access information about the system from your scripts, as well as provide a convenient way to store data within your scripts.

Chapter 5

Using Linux Environment Variables

In This Chapter

- What are environment variables?
- Setting environment variables
- Removing environment variables
- Default shell environment variables
- Setting the PATH environment variable
- Locating system environment variables
- Variable arrays
- Using command aliases

Linux environment variables help define your Linux shell experience. However, they can be a confusing topic for new Linux users. Many programs and scripts use environment variables to obtain system information and store temporary data and configuration information. There are lots of places where environment variables are set on the Linux system, and it's important to know where these places are. This chapter walks you through the world of Linux environment

variables, showing where they are, how to use them, and even how to create your own. The chapter finishes off with a related topic, defining and using aliases in your shell session.

What Are Environment Variables?

The bash shell uses a feature called *environment variables* to store information about the shell session and the working environment (thus the name environment variables). This feature also allows you to store data in memory that can be easily accessed by any program or script running from the shell. This is a handy way to store persistent data that identifies features of the user account, system, shell, or anything else you need to store.

There are two types of environment variables in the bash shell:

- Global variables
- Local variables

This section describes each type of environment variable, and shows how to see and use them.

Note

Even though the bash shell uses specific environment variables that are consistent, different

Linux distributions often add their own environment variables. The environment variable examples you see in this chapter may differ slightly from what's available in your specific distribution. If you run into an environment variable not covered here, check the documentation for your Linux distribution.

Global Environment Variables

Global environment variables are visible from the shell session, and from any child processes that the shell spawns. Local variables are only available in the shell that creates them. This makes global environment variables useful in applications that spawn child processes that require information from the parent process.

The Linux system sets several global environment variables when you start your bash session (for more details about what variables are started at that time, see the “Locating System Environment Variables” section later in this chapter). The system environment variables always use all capital letters to differentiate them from normal user environment variables.

To view the global environment variables, use the `printenv` command:

```
$ printenv
ORBIT_SOCKETDIR=/tmp/orbit-user
HOSTNAME=localhost.localdomain
IMSETTINGS_INTEGRATE_DESKTOP=yes
TERM=xterm
SHELL=/bin/bash
XDG_SESSION_COOKIE=787b3cf537971ef846226096000006b-1284670942.440386-1012435051
HISTSIZE=1000
GTK_RC_FILES=/etc/gtk/gtkrc:/home/user/.gtkrc-1.2-gnome2
WINDOWID=29360131
```

```
GNOME_KEYRING_CONTROL=/tmp/keyring-Egjauk
IMSETTINGS_MODULE=none
USER=user
LS_COLORS=rs=0:di=01;34:ln=01;36:mh=00:pi=40;33:so=01;35:do=01;35:
bd=40;33;01:cd=40;33;01:or=40;31;01:mi=01;05;37;41:su=37;41:sg=30;43
...
SSH_AUTH_SOCK=/tmp/keyring-Egjauk/ssh
SESSION_MANAGER=local/unix:@/tmp/.ICE-unix/1331,unix/unix:/tmp/.ICE-unix/1331
USERNAME=user
DESKTOP_SESSION=gnome
MAIL=/var/spool/mail/user
PATH=/usr/kerberos/sbin:/usr/kerberos/bin:/usr/local/bin:/usr/bin:
/bin:/usr/local/sbin:/usr/sbin:/sbin:/home/user/bin
QT_IM_MODULE=xim
PWD=/home/user/Documents
XMODIFIERS=@im=none
GDM_KEYBOARD_LAYOUT=us
LANG=en_US.utf8
GNOME_KEYRING_PID=1324
GDM_LANG=en_US.utf8
GDMSESSION=gnome
SSH_ASKPASS=/usr/libexec/openssh/gnome-ssh-askpass
HISTCONTROL=ignoredups
HOME=/home/user
SHLVL=2
GNOME_DESKTOP_SESSION_ID=this-is-deprecated
LOGNAME=user
DBUS_SESSION_BUS_ADDRESS=unix:abstract=/tmp/dbus-
ifKtHnDKnU,guid=6f061c07be822c134f12956b00000081
LESSOPEN=|/usr/bin/lesspipe.sh %s
WINDOWPATH=1
DISPLAY=:0.0
G_BROKEN_FILENAMES=1
XAUTHORITY=/var/run/gdm/auth-for-user-GIU7sE/database
COLORTERM=gnome-terminal
_=/usr/bin/printenv
OLDPWD=/home/user
$
```

As you can see, there are lots of global environment variables that get set for the bash shell. Most of them are set by the system during the login process.

To display the value of an individual environment variable, use the `echo` command. When referencing an environment variable, you must place a dollar sign before

the environment variable name:

```
$ echo $HOME
/home/user
$
```

As mentioned, global environment variables are also available to child processes running under the current shell session:

```
$ bash
$ echo $HOME
/home/user
$
```

In this example, after starting a new shell using the `bash` command, you displayed the current value of the `HOME` environment variable, which the system sets when you log into the main shell. Sure enough, the value is also available from the child shell process.

Local Environment Variables

Local environment variables, as their name implies, can be seen only in the local process in which they are defined. Even though they are local, they are just as important as global environment variables. In fact, the Linux system also defines standard local environment variables for you by default.

Trying to see the list of local environment variables is a little tricky. Unfortunately there isn't a command that displays only local environment variables. The `set` command displays all of the environment variables set for a specific process. However, this also includes the global

environment variables.

Here's the output from a sample `set` command:

```
$ set
BASH=/bin/bash
BASHOPTS=checkwinsize:cmdhist:expand_aliases:extquote:
force_ignores:hostcomplete:interactive_comments:
progcomp:promptvars:sourcpath
BASH_ALIASES=()
BASH_ARGC=()
BASH_ARGV=()
BASH_CMDS=()
BASH_LINENO=()
BASH_SOURCE=()
BASH_VERSIONINFO={ [0]="4" [1]="1" [2]="7" [3]="1"
[4]="release" [5]="i386-redhat-linux-gnu" }
BASH_VERSION='4.1.7(1)-release'
COLORS=/etc/DIR_COLORS
COLORTERM=gnome-terminal
COLUMNS=80
DBUS_SESSION_BUS_ADDRESS=unix:abstract=/tmp/dbus-
ifKtHnDKnU,guid=6f061c07be822c134f12956b00000081
DESKTOP_SESSION=gnome
DIRSTACK=()
DISPLAY=:0.0
EUID=500
GDMSESSION=gnome
GDM_KEYBOARD_LAYOUT=us
GDM_LANG=en_US.utf8
GNOME_DESKTOP_SESSION_ID=this-is-deprecated
GNOME_KEYRING_CONTROL=/tmp/keyring-Egjauk
GNOME_KEYRING_PID=1324
GROUPS=()
GTK_RC_FILES=/etc/gtk/gtkrc:/home/user/.gtkrc-1.2-gnome2
G_BROKEN_FILENAMES=1
HISTCONTROL=ignoredups
HISTFILE=/home/user/.bash_history
HISTFILESIZE=1000
HISTSIZ=1000
HOME=/home/user
HOSTNAME=localhost.localdomain
HOSTTYPE=i386
IFS=$' \t\n'
IMSETTINGS_INTEGRATE_DESKTOP=yes
IMSETTINGS_MODULE=none
LANG=en_US.utf8
LESSOPEN='|/usr/bin/lesspipe.sh %s'
LINES=24
LOGNAME=user
```

```
PS_COLORS='rs=0:di=0;34:ln=01;36:mh=00:pi=40;33:so=01;35:do=01;35:
bd=40;33;01:cd=40;33;01:or=40;31;01:mi=01;05;37;41:su=37;41:sg=30;43
...
```

```
MACHTYPE=i386-redhat-linux-gnu
MAIL=/var/spool/mail/user
MAILCHECK=60
OLDPWD=/home/user
OPTERR=1
OPTIND=1
ORBIT_SOCKETDIR=/tmp/orbit-user
OSTYPE=linux-gnu
PATH=/usr/kerberos/sbin:/usr/kerberos/bin:/usr/local/bin:/usr/bin:
/bin:/usr/local/sbin:/usr/sbin:/sbin:/home/user/bin
PIPESTATUS=([0]="0")
PPID=1674
PROMPT_COMMAND='echo -ne "\033]0;${USER}@${HOSTNAME%%.*}:"
${PWD/#$HOME/~}"; echo -ne "\007"'
PS1='\u@\h \W]\$ '
PS2='> '
PS4='+ '
PWD=/home/user/Documents
QT_IM_MODULE=xim
SESSION_MANAGER=local/unix:@/tmp/.ICE-unix/1331,unix/unix:/tmp/.ICE-unix/1331
SHELL=/bin/bash
SHELLOPTS=braceexpand:emacs:hashall:histexpand:history:interactive-
comments:monitor
SHLVL=2
SSH_ASKPASS=/usr/libexec/openssh/gnome-ssh-askpass
SSH_AUTH_SOCK=/tmp/keyring-Egjauk/ssh
TERM=xterm
UID=500
USER=user
USERNAME=user
WINDOWID=29360131
WINDOWPATH=1
XAUTHORITY=/var/run/gdm/auth-for-user-GIU7sE/database
XDG_SESSION_COOKIE=787b3cf537971ef846226096000006b-1284670942.440386-1012435051
XMODIFIERS=@im=none
_=printenv
colors=/etc/DIR_COLORS
__udisks ()
{
    local IFS='
';
    local cur="${COMP_WORDS[COMP_CWORD]}";
    if [ "${COMP_WORDS[${COMP_CWORD} - 1]}" = "--show-info" ]; then
        COMPREPLY=($(compgen -W "${udisks --enumerate-device-files}" -- $cur));
    else
        if [ "${COMP_WORDS[${COMP_CWORD} - 1]}" = "--inhibit-polling" ]; then
```

```

...
        fi;
    fi;
fi;

fi
}

command_not_found_handle ()
{
    runcnf=1;
    retval=127;
    [ ! -S /var/run/dbus/system_bus_socket ] && runcnf=0;
    [ ! -x /usr/sbin/packagekitd ] && runcnf=0;
    if [ $runcnf -eq 1 ]; then
        /usr/libexec/pk-command-not-found $1;
        retval=$?;
    else
        echo "bash: $1: command not found";
    fi;
    return $retval
}

$

```

You'll notice that all of the global environment variables seen from the `printenv` command appear in the output from the `set` command. However, there are quite a few additional environment variables that now appear. These are the local environment variables.

Setting Environment Variables

You can set your own environment variables directly from the bash shell. This section shows you how to create your own environment variables and reference them from your interactive shell or shell script program.

Setting Local Environment Variables

Once you start a bash shell (or spawn a shell script), you're allowed to create local variables that are visible within your shell process. You can assign either a numeric or a string value to an environment variable by assigning the variable to a value using the equal sign:

```
$ test=testing
$ echo $test
testing
$
```

That was simple! Now any time you need to reference the value of the test environment variable, just reference it by the name `$test`.

If you need to assign a string value that contains spaces, you'll need to use a single quotation mark to delineate the beginning and the end of the string:

```
$ test=testing a long string
-bash: a: command not found
$ test='testing a long string'
$ echo $test
testing a long string
$
```

Without the single quotation marks, the bash shell assumes that the next character is another command to process. Notice that for the local environment variable you defined, you used lowercase letters, while the system environment variables you've seen so far have all used uppercase letters.

This is a standard convention in the bash shell. If you create new environment variables, it is recommended (but not required) that you use lowercase letters. This helps distinguish your personal environment variables from the

scores of system environment variables.

Caution

It's extremely important that there are no spaces between the environment variable name, the equal sign, and the value. If you put any spaces in the assignment, the bash shell interprets the value as a separate command:

```
$ test2 = test
-bash: test2: command not found
$
```

Once you set a local environment variable, it's available for use anywhere within your shell process. However, if you spawn another shell, it's not available in the child shell:

```
$ bash
$ echo $test

$ exit
exit
$ echo $test
testing a long string
$
```

In this example, you started a child shell. As you can see, the test environment variable is not available in the child shell (it contains a blank value). After you exited the child shell and returned to the original shell, the local environment variable was still available.

Similarly, if you set a local environment variable in a child process, once you leave the child process the local environment variable is no longer available:

```
$ bash
$ test=testing
$ echo $test
```

```
testing
$ exit
exit
$ echo $test

$
```

The test environment variable set in the child shell doesn't exist when we go back to the parent shell.

Setting Global Environment Variables

Global environment variables are visible from any child processes created by the process that sets the global environment variable. The method used to create a global environment variable is to create a local environment variable and then export it to the global environment.

This is done by using the `export` command:

```
$ echo $test
testing a long string
$ export test
$ bash
$ echo $test
testing a long string
$
```

After exporting the local environment variable `test`, we started a child shell process and viewed the value of the test environment variable. This time, the environment variable kept its value, as the `export` command made it global.

Caution

Notice that when exporting a local environment variable, you

don't use the dollar sign to reference the variable's name.

Removing Environment Variables

Of course, if you can create a new environment variable, it makes sense that you can also remove an existing environment variable. You can do this with the `unset` command:

```
$ echo $test
testing
$ unset test
$ echo $test

$
```

When referencing the environment variable in the `unset` command, remember not to use the dollar sign.

When dealing with global environment variables, things get a little tricky. If you're in a child process and unset a global environment variable, it applies only to the child process. The global environment variable is still available in the parent process:

```
$ test=testing
$ export test
$ bash
$ echo $test
testing
$ unset test
$ echo $test

$ exit
exit
$ echo $test
```

In this example, you set a local environment variable called `test`, and then exported it to make it a global environment variable. You then started a child shell process and checked to make sure that the global environment variable `test` was still available. Next, while still in the child shell, you used the `unset` command to remove the global environment variable `test`, and then exited the child shell. Now back in the original parent shell, you checked the `test` environment variable value, and it is still valid.

Default Shell Environment Variables

There are specific environment variables that the `bash` shell uses by default to define the system environment. You can always count on these variables being set on your Linux system. Because the `bash` shell is a derivative of the original Unix Bourne shell, it also includes environment variables originally defined in that shell.

[Table 5.1](#) shows the environment variables the `bash` shell provides that are compatible with the original Unix Bourne shell.

[Table 5.1](#) The `bash` Shell Bourne Variables

Variable	Description
	A colon-separated list of directories used as a search path for

CDPATH	the <code>cd</code> command.
HOME	The current user's home directory.
IFS	A list of characters that separate fields used by the shell to split text strings.
MAIL	The file name for the current user's mailbox. The bash shell checks this file for new mail.
MAILPATH	A colon-separated list of multiple file names for the current user's mailbox. The bash shell checks each file in this list for new mail.
OPTARG	The value of the last option argument processed by the <code>getopts</code> command.
OPTIND	The index value of the last option argument processed by the <code>getopts</code> command.
PATH	A colon-separated list of directories where the shell looks for commands.
PS1	The primary shell command line interface prompt string.
PS2	The secondary shell command line interface prompt string.

By far the most important environment variable in this list is the `PATH` environment variable. When you enter a command in the shell command line interface (CLI), the shell must search the system to find the program. The `PATH` environment variable defines the directories it searches looking for commands. On my Linux system, the `PATH` environment variable looks like this:

```
$ echo $PATH
/usr/kerberos/sbin:/usr/kerberos/bin:/usr/local/bin:/usr/bin:
/bin:/usr/local/sbin:/usr/sbin:/sbin:/home/user/bin

$
```

This shows that there are six directories where the shell looks for commands. Each directory in the `PATH` is separated

by a colon. There's nothing at the end of the `PATH` variable indicating the end of the directory listing. You can add additional directories to the `PATH` simply by adding another colon, and adding the new directory. The `PATH` also shows the order in which it looks for commands.

Besides the default Bourne environment variables, the bash shell also provides a few variables of its own, as shown in [Table 5.2](#).

Table 5.2 The bash Shell Environment Variables

Variable	Description
<code>BASH</code>	The full pathname to execute the current instance of the bash shell.
<code>BASH_ALIASES</code>	An associative array of currently set aliases.
<code>BASH_ARGC</code>	A variable array that contains the number of parameters being passed to a subroutine or shell script.
<code>BASH_ARCV</code>	A variable array that contains the parameters being passed to a subroutine or shell script.
<code>BASH_CMDS</code>	An associative array of locations of commands the shell has executed.
<code>BASH_COMMAND</code>	The shell command currently being or about to be executed.
<code>BASH_ENV</code>	When set, each bash script attempts to execute a startup file defined by this variable before running.
<code>BASH_EXECUTION_STRING</code>	The command(s) passed using the <code>bash -c</code> option.
<code>BASH_LINENO</code>	A variable array containing the source code line number of the currently executing shell function.
<code>BASH_REMATCH</code>	A read-only variable array containing patterns and their subpatterns for positive matches using the

	regular expression comparison operator, ==.
BASH_SOURCE	A variable array containing the source code file name of the currently executing shell function.
BASH_SUBSHELL	The current nesting level of a subshell environment. The initial value is 0.
BASH_VERSION	The version number of the current instance of the bash shell.
BASH_VERSINFO	A variable array that contains the individual major and minor version numbers of the current instance of the bash shell.
BASH_XTRACEFD	If set to a valid file descriptor (0,1,2), then trace output generated from the 'set -x' debugging option can be redirected. This is often used to separate trace output into a file.
BASHOPTS	A list of bash shell options that are currently enabled.
BASHPID	Process ID of the current bash process.
COLUMNS	Contains the terminal width of the terminal used for the current instance of the bash shell.
COMP_CWORD	An index into the variable COMP_WORDS, which contains the current cursor position.
COMP_LINE	The current command line.
COMP_POINT	The index of the current cursor position relative to the beginning of the current command.
COMP_KEY	The final key used to invoke the current completion of a shell function.
COMP_TYPE	An integer value representing the type of completion attempted that caused a completion shell function to be invoked.
COMP_WORDBREAKS	The Readline library word separator characters for performing word completion.
	A variable array that contains the individual words on

COMP_WORDS	the current command line.
COMPREPLY	A variable array that contains the possible completion codes generated by a shell function.
DIRSTACK	A variable array that contains the current contents of the directory stack.
EMACS	Indicates the emacs shell buffer is executing and line editing is disabled, when set to 't'.
EUID	The numeric effective user ID of the current user.
FCEDIT	The default editor used by the <code>fc</code> command.
FIGIGNORE	A colon-separated list of suffixes to ignore when performing file name completion.
FUNCNAME	The name of the currently executing shell function.
GLOBIGNORE	A colon-separated list of patterns defining the set of file names to be ignored by file name expansion.
GROUPS	A variable array containing the list of groups of which the current user is a member.
histchars	Up to three characters, which control history expansion.
HISTCMD	The history number of the current command.
HISTCONTROL	Controls what commands are entered in the shell history list.
HISTFILE	The name of the file to save the shell history list (<code>.bash_history</code> by default).
HISTFILESIZE	The maximum number of lines to save in the history file.
HISTIGNORE	A colon-separated list of patterns used to decide which commands are ignored for the history file.
HISTSIZE	The maximum number of commands stored in the history file.
	Contains the name of the file that should be read

HOSTFILE	when the shell needs to complete a hostname.
HOSTNAME	The name of the current host.
HOSTTYPE	A string describing the machine the bash shell is running on.
IGNOREEOF	The number of consecutive EOF characters the shell must receive before exiting. If this value doesn't exist, the default is 1.
INPUTRC	The name of the Readline initialization file (the default is <code>.inputrc</code>).
LANG	The locale category for the shell.
LC_ALL	Overrides the <code>LANG</code> variable, defining a locale category.
LC_COLLATE	Sets the collation order used when sorting string values.
LC_CTYPE	Determines the interpretation of characters used in file name expansion and pattern matching.
LC_MESSAGES	Determines the locale setting used when interpreting double-quoted strings preceded by a dollar sign.
LC_NUMERIC	Determines the locale setting used when formatting numbers.
LINENO	The line number in a script currently executing.
LINES	Defines the number of lines available on the terminal.
MACHTYPE	A string defining the system type in <code>cpu-company-system</code> format.
MAILCHECK	How often (in seconds) the shell should check for new mail (default is 60).
OLDPWD	The previous working directory used in the shell.
OPTERR	If set to 1, the bash shell displays errors generated by the <code>getopts</code> command.

OSTYPE	A string defining the operating system the shell is running on.
PIPESTATUS	A variable array containing a list of exit status values from the processes in the foreground process.
POSIXLY_CORRECT	If set, bash starts in POSIX mode.
PPID	The process ID (PID) of the bash shell's parent process.
PROMPT_COMMAND	If set, the command to execute before displaying the primary prompt.
PROMPT_DIRTRIM	An integer used to indicate the number of trailing directory names to display when using the \w and \W prompt string escapes. The directory names removed are replaced with one set of ellipses.
PS3	The prompt to use for the <code>select</code> command.
PS4	The prompt displayed before the command line is echoed if the bash <code>-x</code> parameter is used.
PWD	The current working directory.
RANDOM	Returns a random number between 0 and 32767. Assigning a value to this variable seeds the random number generator.
REPLY	The default variable for the <code>read</code> command.
SECONDS	The number of seconds since the shell was started. Assigning a value resets the timer to the value.
SHELL	The full path name to the bash shell.
SHELLOPTS	A colon-separated list of enabled bash shell options.
SHLVL	Indicates the shell level, incremented by one each time a new bash shell is started.
TIMEFORMAT	A format specifying how the shell displays time values.
TMOUT	The value of how long (in seconds) the <code>select</code> and <code>read</code> commands should wait for input. The default of

	zero indicates to wait indefinitely.
<code>TMPDIR</code>	Directory name where the bash shell creates temporary files for its use.
<code>UID</code>	The numeric real user ID of the current user.

You may notice that not all of the default environment variables are shown when the `set` command was used. The reason for this is that although these are the default environment variables, not all of them are required to contain a value.

Setting the PATH Environment Variable

The `PATH` environment variable seems to cause the most problem on Linux systems. It defines where the shell looks for commands you enter on the command line. If it can't find the command, it produces an error message:

```
$ myprog
-bash: myprog: command not found
$
```

The problem is that often applications place their executable programs in directories that aren't in the `PATH` environment variable. The trick is ensuring that your `PATH` environment variable includes all of the directories where your applications reside.

You can add new search directories to the existing `PATH` environment variable without having to rebuild it from

scratch. The individual directories listed in the `PATH` are separated by a colon. All you need to do is reference the original `PATH` value and add any new directories to the string. This looks something like this:

```
$ echo $PATH
/usr/kerberos/sbin:/usr/kerberos/bin:/usr/local/bin:/usr/bin:
/bin:/usr/local/sbin:/usr/sbin:/sbin:/home/user/bin

$ PATH=$PATH:/home/user/test
$ echo $PATH
/usr/kerberos/sbin:/usr/kerberos/bin:/usr/local/bin:/usr/bin:
/bin:/usr/local/sbin:/usr/sbin:/sbin:/home/user/bin:/home/user/test
$ myprog
The factorial of 5 is 120.
$
```

By adding the directory to the `PATH` environment variable, you can now execute your program from anywhere in the virtual directory structure:

```
[user@localhost ~]$ cd /etc
[user@localhost etc]$ myprog
The factorial of 5 is 120
[user@localhost etc]$
```

A common trick for programmers is to include the single dot symbol in their `PATH` environment variable. The single dot symbol represents the current directory (see Chapter 3):

```
[user@localhost ~]$ PATH=$PATH:
[user@localhost ~]$ cd test2
[user@localhost test2]$ myprog2
The factorial of 6 is 720
[user@localhost test2]$
```

In the next section, you'll see how you can make changes to environment variables permanent on your system, so you can always execute your programs.

Locating System Environment Variables

The Linux system uses environment variables to identify itself in programs and scripts. This provides a convenient way to obtain system information for your programs. The trick is in how these environment variables are set.

When you start a bash shell by logging in to the Linux system, by default bash checks several files for commands. These files are called *startup files*. The startup files bash processes depend on the method you use to start the bash shell. There are three ways of starting a bash shell:

- As a default login shell at login time
- As an interactive shell that is not the login shell
- As a non-interactive shell to run a script

The following sections describe the startup files the bash shell executes in each of these startup methods.

Login Shell

When you log in to the Linux system, the bash shell starts as a login shell. The login shell looks for four different startup files to process commands from. The following is the order in which the bash shell processes the files:

- `/etc/profile`
- `$HOME/.bash_profile`
- `$HOME/.bash_login`

- `$HOME/.profile`

The `/etc/profile` file is the main default startup file for the bash shell on the system. Every user on the system executes this startup file when they log in. The other three startup files are specific for each user and can be customized for each user's requirements. Let's take a closer look at these files.

The /etc/profile File

The `/etc/profile` file is the main default startup file for the bash shell. Whenever you log in to the Linux system, bash executes the commands in the `/etc/profile` startup file. Different Linux distributions place different commands in this file. On this Linux system, it looks like this:

```
$ cat /etc/profile
# /etc/profile

# System wide environment and startup programs, for login setup
# Functions and aliases go in /etc/bashrc

# It's NOT good idea to change this file unless you know what you
# are doing. Much better way is to create custom.sh shell script in
# /etc/profile.d/ to make custom changes to environment. This will
# prevent need for merging in future updates.

pathmunge () {
    case "${PATH}" in
        *:"$1":*)
            ;;
        *)
            if [ "$2" = "after" ] ; then
                PATH=$PATH:$1
            else
                PATH=$1:$PATH
            fi
    esac
}
```

```

if [ -x /usr/bin/id ]; then
    if [ -z "$EUID" ]; then
        # ksh workaround
        EUID='id -u'
        UID='id -ru'
    fi
    USER="'id -un'"
    LOGNAME=$USER
    MAIL="/var/spool/mail/$USER"
fi

# Path manipulation
if [ "$EUID" = "0" ]; then
    pathmunge /sbin
    pathmunge /usr/sbin
    pathmunge /usr/local/sbin
else
    pathmunge /usr/local/sbin after
    pathmunge /usr/sbin after
    pathmunge /sbin after
fi

HOSTNAME='/bin/hostname 2>/dev/null'
HISTSIZE=1000
if [ "$HISTCONTROL" = "ignorespace" ] ; then
    export HISTCONTROL=ignoreboth
else
    export HISTCONTROL=ignoredups
fi

export PATH USER LOGNAME MAIL HOSTNAME HISTSIZE HISTCONTROL

for i in /etc/profile.d/*.sh ; do
    if [ -r "$i" ]; then
        if [ "$PS1" ]; then
            . $i
        else
            . $i >/dev/null 2>&1
        fi
    fi
done

unset i
unset pathmunge
$

```

Most of the commands and scripts you see in this file are covered in more detail later on in Chapter 10. The

important thing to notice now is the environment variables that are set in this startup file. Notice the export line near the bottom of the file:

```
export PATH USER LOGNAME MAIL HOSTNAME HISTSIZE HISTCONTROL
```

This ensures that these environment variables are available to all child processes spawned from the login shell.

There's also another tricky feature that the profile file uses. There's a `for` statement that iterates through any files located in the `/etc/profile.d` directory. (`for` statements are discussed in detail in Chapter 12.) This provides a place for the Linux system to place application-specific startup files that will be executed by the shell when you log in. On this Linux system, the following files are in the `profile.d` directory:

```
$ ls -l /etc/profile.d
total 72
-rw-r--r--. 1 root root 1133 Apr 28 10:43 colorls.csh
-rw-r--r--. 1 root root 1143 Apr 28 10:43 colorls.sh
-rw-r--r--. 1 root root 192 Sep 9 2004 glib2.csh
-rw-r--r--. 1 root root 192 Dec 12 2005 glib2.sh
-rw-r--r--. 1 root root 58 May 31 06:23 gnome-ssh-askpass.csh
-rw-r--r--. 1 root root 70 May 31 06:23 gnome-ssh-askpass.sh
-rw-r--r--. 1 root root 184 Aug 25 11:36 krb5-workstation.csh
-rw-r--r--. 1 root root 57 Aug 25 11:36 krb5-workstation.sh
-rw-r--r--. 1 root root 1741 Jun 24 15:20 lang.csh
-rw-r--r--. 1 root root 2706 Jun 24 15:20 lang.sh
-rw-r--r--. 1 root root 122 Feb 7 2007 less.csh
-rw-r--r--. 1 root root 108 Feb 7 2007 less.sh
-rw-r--r--. 1 root root 837 Sep 2 05:24 PackageKit.sh
-rw-r--r--. 1 root root 2142 Aug 10 16:41 udisks-bash-completion.sh
-rw-r--r--. 1 root root 74 Mar 25 19:24 vim.csh
-rw-r--r--. 1 root root 248 Mar 25 19:24 vim.sh
-rw-r--r--. 1 root root 161 Nov 27 2007 which2.csh
-rw-r--r--. 1 root root 169 Nov 27 2007 which2.sh
$
```

You'll notice that these are mostly related to specific

applications on the system. Most applications create two startup files, one for the bash shell (using the `.sh` extension) and one for the c shell (using the `.csh` extension).

The `lang.csh` and `lang.sh` files attempt to determine the default language character set used on the system, and set the `LANG` environment variable appropriately.

The \$HOME Startup Files

The remaining three startup files are all used for the same function—to provide a user-specific startup file for defining user-specific environment variables. Most Linux distributions use only one of these three startup files:

- `$HOME/.bash_profile`
- `$HOME/.bash_login`
- `$HOME/.profile`

Notice that all three files start with a dot, making them hidden files (they don't appear in a normal `ls` command listing). Because they are in the user's `HOME` directory, each user can edit the files and add his or her own environment variables that are active for every bash shell session they start.

This Linux system contains the following `.bash_profile` file:

```
$ cat .bash_profile
# .bash_profile

# Get the aliases and functions
if [ -f ~/.bashrc ]; then
    . ~/.bashrc
fi

# User specific environment and startup programs
```

```
PATH=$PATH:$HOME/bin
```

```
export PATH
```

```
$
```

The `.bash_profile` startup file first checks to see if there's another startup file present in the `HOME` directory, called `.bashrc` (which we'll talk about next in the “Interactive Shell” section). If it's there, the startup file executes the commands in it. Next, the startup file adds a directory to the `PATH` environment variable, providing a common location to place executable files in your `HOME` directory.

Interactive Shell

If you start a bash shell without logging into a system (such as if you just type `bash` at a CLI prompt), you start what's called an *interactive shell*. The interactive shell doesn't act like the login shell, but it still provides a CLI prompt for you to enter commands.

If `bash` is started as an interactive shell, it doesn't process the `/etc/profile` file. Instead, it checks for the `.bashrc` file in the user's `HOME` directory.

On this Linux distribution, this file looks like this:

```
$ cat .bashrc
```

```
# .bashrc
```

```
# Source global definitions
```

```
if [ -f /etc/bashrc ]; then
```

```
    . /etc/bashrc
```

```
fi
```

```
# User specific aliases and functions
```

```
$
```

The `.bashrc` file does two things. First, it checks for a common `bashrc` file in the `/etc` directory. Second, it provides a place for the user to enter personal aliases (discussed later in the “Using Command Aliases” section) and private script functions (described in Chapter 16).

The common `/etc/bashrc` startup file is run by everyone on the system who starts an interactive shell session. On this Linux distribution it looks like this:

```
$ cat /etc/bashrc
# /etc/bashrc

# System wide functions and aliases
# Environment stuff goes in /etc/profile

# It's NOT good idea to change this file unless you know what you
# are doing. Much better way is to create custom.sh shell script in
# /etc/profile.d/ to make custom changes to environment. This will
# prevent need for merging in future updates.

# By default, we want this to get set.
# Even for non-interactive, non-login shells.
# Current threshold for system reserved uid/gids is 200
# You could check uidgid reservation validity in
# /usr/share/doc/setup-*/uidgid file
if [ $UID -gt 199 ] && [ "`id -gn`" = "`id -un`" ]; then
    umask 002
else
    umask 022
fi

# are we an interactive shell?
if [ "$PS1" ]; then
    case $TERM in
        xterm*)
            if [ -e /etc/sysconfig/bash-prompt-xterm ]; then
                PROMPT_COMMAND=/etc/sysconfig/bash-prompt-xterm
            else
                PROMPT_COMMAND='echo -ne
"\033]0;${USER}@${HOSTNAME%%.*}:
${PWD/#$HOME/~}"; echo -ne "\007"'
            fi
        ;;
        screen)
            if [ -e /etc/sysconfig/bash-prompt-screen ]; then
```

```

        PROMPT_COMMAND=/etc/sysconfig/bash-prompt-screen
    else
        PROMPT_COMMAND='echo -ne
"\033_${USER}@${HOSTNAME%%.*}:
${PWD/#$HOME/~}"; echo -ne "\033\`"
    fi
;;
*)
    [ -e /etc/sysconfig/bash-prompt-default ] &&
    PROMPT_COMMAND=/etc/sysconfig/bash-prompt-default
;;
esac
# Turn on checkwinsize
shopt -s checkwinsize
[ "$PS1" = "\\s-\\v\\\$ " ] && PS1="[\\u@\\h \\W]\\\$ "
# You might want to have e.g. tty in prompt (e.g. more virtual machines)
# and console windows
# If you want to do so, just add e.g.
# if [ "$PS1" ]; then
#     PS1="[\\u@\\h:\\l \\W]\\\$ "
# fi
# to your custom modification shell script in /etc/profile.d/ directory
fi

if ! shopt -q login_shell ; then # We're not a login shell
    # Need to redefine pathmunge, it get's undefined at the end of /etc/profile
    pathmunge () {
        case ":${PATH}:" in
            *:"$1":*)
                ;;
            *)
                if [ "$2" = "after" ] ; then
                    PATH=$PATH:$1
                else
                    PATH=$1:$PATH
                fi
            esac
        }

    # Only display echos from profile.d scripts if we are no login shell
    # and interactive - otherwise just process them to set envvars
    for i in /etc/profile.d/*.sh; do
        if [ -r "$i" ]; then
            if [ "$PS1" ]; then
                . $i
            else
                . $i >/dev/null 2>&1
            fi
        fi
    done

```



```
unset i
unset pathmunge
fi
# vim:ts=4:sw=4
$
```

The default file sets a few environment variables, but notice that it doesn't use the `export` command to make them global. Remember, the interactive shell startup file runs each time a new interactive shell starts; thus, any child shell will automatically execute the interactive shell startup file.

You'll also notice that the `/etc/bashrc` file also executes the application-specific startup files located in the `/etc/profile.d` directory.

Non-Interactive Shell

Finally, the last type of shell is a non-interactive shell. This is the shell that the system starts to execute a shell script. This is different in that there isn't a CLI prompt to worry about. However, there may still be specific startup commands you want to run each time you start a script on your system.

To accommodate that situation, the bash shell provides the `BASH_ENV` environment variable. When the shell starts a non-interactive shell process, it checks this environment variable for the name of a startup file to execute. If one is present, the shell executes the commands in the file. On our Linux distribution, this environment value is not set by default.

Variable Arrays

A really cool feature of environment variables is that they can be used as *arrays*. An array is a variable that can hold multiple values. Values can be referenced either individually or as a whole for the entire array.

To set multiple values for an environment variable, just list them in parentheses, with each value separated by a space:

```
$ mytest=(one two three four five)
$
```

Not much excitement there. If you try to display the array as a normal environment variable, you'll be disappointed:

```
$ echo $mytest
one
$
```

Only the first value in the array appears. To reference an individual array element, you must use a numerical index value, which represents its place in the array. The numeric value is enclosed in square brackets:

```
$ echo ${mytest[2]}
three
$
```

Caution

Environment variable arrays start with an index value of zero. This often gets confusing.

To display an entire array variable, you use the asterisk

wildcard character as the index value:

```
$ echo ${mytest[*]}  
one two three four five  
$
```

You can also change the value of an individual index position:

```
$ mytest[2]=seven  
$ echo ${mytest[*]}  
one two seven four five  
$
```

You can even use the `unset` command to remove an individual value within the array, but be careful, as this gets tricky. Watch this example:

```
$ unset mytest[2]  
$ echo ${mytest[*]}  
one two four five  
$  
$ echo ${mytest[2]}  
  
$ echo ${mytest[3]}  
four  
$
```

This example uses the `unset` command to remove the value at index value 2. When you display the array, it appears that the other index values just dropped down one. However, if you specifically display the data at index value 2, you'll see that that location is empty.

Finally, you can remove the entire array just by using the array name in the `unset` command:

```
$ unset mytest  
$ echo ${mytest[*]}  
  
$
```

Sometimes variable arrays just complicate matters, so

they're often not used in shell script programming. They're not very portable to other shell environments, which is a downside if you do lots of shell programming for different shells. There are a couple of bash system environment variables that use arrays (such as `BASH_VERSINFO`), but overall you probably won't run into them very often.

Using Command Aliases

While not officially environment variables, shell command aliases behave in much the same manner. A *command alias* allows you to create an alias name for common commands (along with their parameters) to help keep your typing to a minimum.

Most likely, your Linux distribution has already set some common command aliases for you. To see a list of the active aliases, use the `alias` command with the `-p` parameter:

```
$ alias -p
alias l.='ls -d .* --color=auto'
alias ll='ls -l --color=auto'
alias ls='ls --color=auto'
alias vi='vim'
alias which='alias | /usr/bin/which --tty-only --read-alias --show-dot
--show-tilde'
$
```

Notice that on this Linux distribution, an alias is used to override the standard `ls` command. It automatically provides the `--color` parameter, indicating that the terminal supports color mode listings.

You can create your own aliases by using the `alias` command:

```
$ alias li='ls -il'
$ li
total 32
75 drwxr-xr-x. 2 user user 4096 Sep 16 13:11 Desktop
79 drwxr-xr-x. 2 user user 4096 Sep 20 15:40 Documents
76 drwxr-xr-x. 2 user user 4096 Sep 15 13:30 Downloads
80 drwxr-xr-x. 2 user user 4096 Sep 15 13:30 Music
81 drwxr-xr-x. 2 user user 4096 Sep 15 13:30 Pictures
78 drwxr-xr-x. 2 user user 4096 Sep 15 13:30 Public
77 drwxr-xr-x. 2 user user 4096 Sep 15 13:30 Templates
82 drwxr-xr-x. 2 user user 4096 Sep 15 13:30 Videos
$
```

Once you define an alias value, you can use it at any time in your shell, including in shell scripts.

Command aliases act like local environment variables. They're only valid for the shell process in which they're defined:

```
$ alias li='ls -il'
$ bash
$ li
bash: li: command not found
$
```

Of course, now you know a way to solve that problem. The `bash` shell always reads the `$HOME/.bashrc` startup file when starting a new interactive shell. This is a great place to put command `alias` statements (as was pointed out in the `.bashrc` file comments).

Summary

This chapter examined the world of Linux environment variables. Global environment variables can be accessed

from any child process spawned by the process they're defined in. Local environment variables can only be accessed from the process in which they're defined.

The Linux system uses both global and local environment variables to store information about the system environment. You can access this information from the shell command line interface, as well as within shell scripts. The bash shell uses the system environment variables defined in the original Unix Bourne shell, as well as lots of new environment variables. The `PATH` environment variable defines the search pattern the bash shell takes to find an executable command. You can modify the `PATH` environment variable to add your own directories, or even the current directory symbol, to make running your programs easier.

You can also create your own global and local environment variables for your own use. Once you create an environment variable, it's accessible for the entire duration of your shell session.

There are several startup files that the bash shell executes when it starts up. These startup files can contain environment variable definitions to set standard environment variables for each bash session. When you log in to the Linux system, the bash shell accesses the `/etc/profile` startup file, and also three local startup files for each user, `$HOME/.bash_profile`, `$HOME/.bash_login`, and `$HOME/.profile`. Users can customize these files to include environment variables and startup scripts for their own use.

The bash shell also provides for environment variable

arrays. These environment variables can contain multiple values in a single variable. You can access the values either individually by referencing an index value or as a whole by referencing the entire environment variable array name.

Finally, the chapter discussed the use of command aliases. While not environment variables, command aliases behave much like environment variables. They allow you to define an alias name for a command, along with its parameters. Instead of having to type in a long command and parameters, you can just assign it to a simple alias and use the alias at any time in your shell session.

The next chapter dives into the world of Linux file permissions. This is possibly the most difficult topic for novice Linux users. However, to write good shell scripts, you need to understand how file permissions work and be able to use them in your Linux system.

Chapter 6

Understanding Linux File Permissions

In This Chapter

- Understanding Linux security
- Decoding the permissions
- Working with Linux groups

No system is complete without some form of security. There must be a mechanism available to protect files from unauthorized viewing or modification. The Linux system follows the Unix method of file permissions, allowing individual users and groups access to files based on a set of security settings for each file and directory. This chapter discusses how to use the Linux file security system to protect data when necessary and share data when desired.

Linux Security

The core of the Linux security system is the *user account*.

Each individual who accesses a Linux system should have a unique user account assigned. The permissions users have to objects on the system depend on the user account they log in with.

User permissions are tracked using a *user ID* (often called a **UID**), which is assigned to an account when it's created. The UID is a numerical value, unique for each user. However, you don't log in to a Linux system using your UID. Instead, you use a *login name*. The login name is an alphanumeric text string of eight characters or fewer that the user uses to log in to the system (along with an associated password).

The Linux system uses special files and utilities to track and manage user accounts on the system. Before we can discuss file permissions, we need to discuss how Linux handles user accounts. This section describes the files and utilities required for user accounts so that you can understand how to use them when working with file permissions.

The `/etc/passwd` File

The Linux system uses a special file to match the login name to a corresponding UID value. This file is the `/etc/passwd` file. The `/etc/passwd` file contains several pieces of information about the user. Here's what a typical `/etc/passwd` file looks like on a Linux system:

```
$ cat /etc/passwd
root:x:0:0:root:/root:/bin/bash
bin:x:1:1:bin:/bin:/sbin/nologin
```

```
daemon:x:2:2:daemon:/sbin:/sbin/nologin
adm:x:3:4:adm:/var/adm:/sbin/nologin
lp:x:4:7:lp:/var/spool/lpd:/sbin/nologin
sync:x:5:0:sync:/sbin:/bin/sync
shutdown:x:6:0:shutdown:/sbin:/sbin/shutdown
halt:x:7:0:halt:/sbin:/sbin/halt
mail:x:8:12:mail:/var/spool/mail:/sbin/nologin
news:x:9:13:news:/etc/news:
uucp:x:10:14:uucp:/var/spool/uucp:/sbin/nologin
operator:x:11:0:operator:/root:/sbin/nologin
games:x:12:100:games:/usr/games:/sbin/nologin
gopher:x:13:30:gopher:/var/gopher:/sbin/nologin
ftp:x:14:50:FTP User:/var/ftp:/sbin/nologin
nobody:x:99:99:Nobody:/:/sbin/nologin
rpm:x:37:37:/:/var/lib/rpm:/sbin/nologin
vcsa:x:69:69:virtual console memory owner:/dev:/sbin/nologin
mailnull:x:47:47:/:/var/spool/mqueue:/sbin/nologin
smmsp:x:51:51:/:/var/spool/mqueue:/sbin/nologin
apache:x:48:48:Apache:/var/www:/sbin/nologin
rpc:x:32:32:Rpcbind Daemon:/var/lib/rpcbind:/sbin/nologin
ntp:x:38:38:/:etc/ntp:/sbin/nologin
nscd:x:28:28:NSCD Daemon:/:/sbin/nologin
tcpdump:x:72:72:/:/sbin/nologin
dbus:x:81:81:System message bus:/:/sbin/nologin
avahi:x:70:70:Avahi daemon:/:/sbin/nologin
hsqldb:x:96:96:/:var/lib/hsqldb:/sbin/nologin
sshd:x:74:74:Privilege-separated SSH:/var/empty/sshd:/sbin/nologin
rpcuser:x:29:29:RPC Service User:/var/lib/nfs:/sbin/nologin
nfsnobody:x:65534:65534:Anonymous NFS User:/var/lib/nfs:/sbin/nologin
haldaemon:x:68:68:HAL daemon:/:/sbin/nologin
xfs:x:43:43:X Font Server:/etc/X11/fs:/sbin/nologin
gdm:x:42:42:/:var/gdm:/sbin/nologin
rich:x:500:500:Rich Blum:/home/rich:/bin/bash
mama:x:501:501:Mama:/home/mama:/bin/bash
katie:x:502:502:katie:/home/katie:/bin/bash
jessica:x:503:503:Jessica:/home/jessica:/bin/bash
mysql:x:27:27:MySQL Server:/var/lib/mysql:/bin/bash
$
```

The `root` user account is the administrator for the Linux system and is always assigned UID 0. As you can see, the Linux system creates lots of user accounts for various functions that aren't actual users. These are called *system accounts*. A system account is a special account that services running on the system use to gain access to

resources on the system. All services that run in background mode need to be logged in to the Linux system under a system user account.

Before security became a big issue, these services often just logged in using the `root` user account. Unfortunately, if an unauthorized person broke into one of these services, he instantly gained access to the system as the `root` user. To prevent this, now just about every service that runs in background on a Linux server has its own user account to log in with. This way, if a troublemaker does compromise a service, he still can't necessarily get access to the whole system.

Linux reserves UIDs below 500 for system accounts. Some services even require specific UIDs to work properly. When you create accounts for normal users, most Linux systems assign the first available UID starting at 500 (although this is not necessarily true for all Linux distributions).

You probably noticed that the `/etc/passwd` file contains lots more than just the login name and UID for the user. The fields of the `/etc/passwd` file contain the following information:

- The login username
- The password for the user
- The numerical UID of the user account
- The numerical group ID (GID) of the user account
- A text description of the user account (called the comment field)
- The location of the `HOME` directory for the user

- The default shell for the user

The password field in the `/etc/passwd` file is set to an `x`. This doesn't mean that all of the user accounts have the same password. In the old days of Linux, the `/etc/passwd` file contained an encrypted version of the user's password. However, since lots of programs need to access the `/etc/passwd` file for user information, this became somewhat of a security problem. With the advent of software that could easily decrypt encrypted passwords, the bad guys had a field day trying to break user passwords stored in the `/etc/passwd` file. Linux developers needed to rethink that policy.

Now, most Linux systems hold user passwords in a separate file (called the *shadow* file, located at `/etc/shadow`). Only special programs (such as the login program) are allowed access to this file.

As you can see, the `/etc/passwd` file is a standard text file. You can use any text editor to manually perform user management functions (such as adding, modifying, or removing user accounts) directly in the `/etc/passwd` file. However, this is an extremely dangerous practice. If the `/etc/passwd` file becomes corrupt, the system won't be able to read it, and it will prevent anyone (even the root user) from logging in. Instead, it's safer to use the standard Linux user management utilities to perform all user management functions.

The `/etc/shadow` File

The `/etc/shadow` file provides more control over how the Linux system manages passwords. Only the root user has access to the `/etc/shadow` file, making it more secure than the `/etc/passwd` file.

The `/etc/shadow` file contains one record for each user account on the system. A record looks like this:

```
rich:$1$.FfcK0ns$f1UgiyHQ25wrB/hykCn020:11627:0:99999:7:::
```

There are nine fields in each `/etc/shadow` file record:

- The login name corresponding to the login name in the `/etc/passwd` file
- The encrypted password
- The number of days since January 1, 1970, that the password was last changed
- The minimum number of days before the password can be changed
- The number of days before the password must be changed
- The number of days before password expiration that the user is warned to change the password
- The number of days after a password expires before the account will be disabled
- The date (stored as the number of days since January 1, 1970) since the user account was disabled
- A field reserved for future use

Using the shadow password system, the Linux system has much finer control over user passwords. It can control

how often a user must change his or her password, and when to disable the account if the password hasn't been changed.

Adding a New User

The primary tool used to add new users to your Linux system is `useradd`. This command provides an easy way to create a new user account and set up the user's `HOME` directory structure all at once. The `useradd` command uses a combination of system default values and command line parameters to define a user account. To see the system default values used on your Linux distribution, enter the `useradd` command with the `-D` parameter:

```
# /usr/sbin/useradd -D
GROUP=100
HOME=/home
INACTIVE=-1
EXPIRE=
SHELL=/bin/bash
SKEL=/etc/skel
CREATE_MAIL_SPOOL=yes
#
```

Note

Some Linux distributions place the Linux user and group utilities in the `/usr/sbin` directory, which may not be in your `PATH` environment variable. If that's the case in your Linux distribution, either add the directory to your `PATH` or use the absolute filepath to run it.

The `-D` parameter shows what defaults the `useradd`

command uses if you don't specify them in the command line when creating a new user account. This example shows the following default values:

- The new user will be added to a common group with group ID 100.
- The new user will have a `HOME` account created in the directory `/home/loginname`.
- The account will not be disabled when the password expires.
- The new account will not be set to expire at a set date.
- The new account will use the bash shell as the default shell.
- The system will copy the contents of the `/etc/skel` directory to the user's `HOME` directory.
- The system will create a file in the mail directory for the user account to receive mail.

The penultimate value is interesting. The `useradd` command allows an administrator to create a default `HOME` directory configuration and then uses that as a template to create the new user's `HOME` directory. This allows you to place default files for the system in every new user's `HOME` directory automatically. In the Ubuntu Linux system, the `/etc/skel` directory has the following files:

```
$ ls -al /etc/skel
total 32
drwxr-xr-x  2 root root  4096 2010-04-29 08:26 .
drwxr-xr-x 135 root root 12288 2010-09-23 18:49 ..
-rw-r--r--  1 root root   220 2010-04-18 21:51 .bash_logout
-rw-r--r--  1 root root  3103 2010-04-18 21:51 .bashrc
```

```
-rw-r--r-- 1 root root 179 2010-03-26 08:31 examples.desktop
-rw-r--r-- 1 root root 675 2010-04-18 21:51 .profile
$
```

You should recognize these files from Chapter 5. These are the standard startup files for the bash shell environment. The system automatically copies these default files into every user's `HOME` directory you create.

You can test this by creating a new user account using the default system parameters and then looking at the `HOME` directory for the new user:

```
# useradd -m test
# ls -al /home/test
total 24
drwxr-xr-x 2 test test 4096 2010-09-23 19:01 .
drwxr-xr-x 4 root root 4096 2010-09-23 19:01 ..
-rw-r--r-- 1 test test 220 2010-04-18 21:51 .bash_logout
-rw-r--r-- 1 test test 3103 2010-04-18 21:51 .bashrc
-rw-r--r-- 1 test test 179 2010-03-26 08:31 examples.desktop
-rw-r--r-- 1 test test 675 2010-04-18 21:51 .profile
#
```

By default, the `useradd` command doesn't create a `HOME` directory, but the `-m` command line option tells it to create the `HOME` directory. As you can see in the example, the `useradd` command created the new `HOME` directory, using the files contained in the `/etc/skel` directory.

Note

To run the user account administration commands in this chapter, you either need to be logged in as the special root user account or use the `sudo` command to run the commands as the root user account.

If you want to override a default value or behavior when creating a new user, you can do that with command line parameters. These are shown in [Table 6.1](#).

Table 6.1 The useradd Command Line Parameters

Parameter	Description
-c <i>comment</i>	Add text to the new user's comment field.
-d <i>home_dir</i>	Specify a different name for the home directory other than the login name.
-e <i>expire_date</i>	Specify a date, in YYYY-MM-DD format, when the account will expire.
-f <i>inactive_days</i>	Specify the number of days after a password expires when the account will be disabled. A value of 0 disables the account as soon as the password expires; a value of -1 disables this feature.
-g <i>initial_group</i>	Specify the group name or GID of the user's login group.
-G <i>group . . .</i>	Specify one or more supplementary groups the user belongs to.
-k	Copy the <code>/etc/skel</code> directory contents into the user's <code>HOME</code> directory (must use <code>-m</code> as well).
-m	Create the user's <code>HOME</code> directory.
-M	Don't create a user's <code>HOME</code> directory (used if the default setting is to create one).
-n	Create a new group using the same name as the user's login name.
-r	Create a system account
-p <i>passwd</i>	Specify a default password for the user account.
-s <i>shell</i>	Specify the default login shell.
-u <i>uid</i>	Specify a unique UID for the account.

As you can see, you can override all of the system default values when creating a new user account just by using command line parameters. However, if you find yourself having to override a value all the time, it's easier to just change the system default value.

You can change the system default new user values by using the `-D` parameter, along with a parameter representing the value you need to change. These parameters are shown in [Table 6.2](#).

Table 6.2 The useradd Change Default Values Parameters

Parameter	Description
<code>-b default_home</code>	Change the location where users' HOME directories are created.
<code>-e expiration_date</code>	Change the expiration date on new accounts.
<code>-f inactive</code>	Change the number of days after a password has expired before the account is disabled.
<code>-g group</code>	Change the default group name or GID used.
<code>-s shell</code>	Change the default login shell.

Changing the default values is a snap:

```
# useradd -D -s /bin/ttsch
# useradd -D
GROUP=100
HOME=/home
INACTIVE=-1
EXPIRE=
SHELL=/bin/ttsch
SKEL=/etc/skel
CREATE_MAIL_SPOOL=yes
#
```

Now, the `useradd` command will use the `tsch` shell as the default login shell for all new user accounts you create.

Removing a User

If you want to remove a user from the system, the `userdel` command is what you need. By default, the `userdel` command removes only the user information from the `/etc/passwd` file. It doesn't remove any files the account owns on the system.

If you use the `-r` parameter, `userdel` will remove the user's `HOME` directory, along with the user's mail directory. However, there may still be other files owned by the deleted user account on the system. This can be a problem in some environments.

Here's an example of using the `userdel` command to remove an existing user account:

```
# /usr/sbin/userdel -r test
# ls -al /home/test
ls: cannot access /home/test: No such file or directory
#
```

After using the `-r` parameter, the user's old `/home/test` directory no longer exists.

Caution

Be careful when using the `-r` parameter in an environment with lots of users. You never know if a user had important files stored in his or her `HOME` directory that are used by someone else or another program. Always check before removing a user's `HOME` directory!

Modifying a User

Linux provides a few different utilities for modifying the information for existing user accounts. [Table 6.3](#) shows these utilities.

Table 6.3 User Account Modification Utilities

Command	Description
<code>usermod</code>	Edits user account fields, as well as specifying primary and secondary group membership
<code>passwd</code>	Changes the password for an existing user
<code>chpasswd</code>	Reads a file of login name and password pairs, and updates the passwords
<code>chage</code>	Changes the password's expiration date
<code>chfn</code>	Changes the user account's comment information
<code>chsh</code>	Changes the user account's default shell

Each of these utilities provides a specific function for changing information about user accounts. The following sections describe each of these utilities.

usermod

The `usermod` command is the most robust of the user account modification utilities. It provides options for changing most of the fields in the `/etc/passwd` file. To do that, you just need to use the command line parameter that corresponds to the value you want to change. The parameters are mostly the

same as the `useradd` parameters (such as `-c` to change the comment field, `-e` to change the expiration date, and `-g` to change the default login group). However, there are a couple of additional parameters that might come in handy:

- `-l` to change the login name of the user account
- `-L` to lock the account so the user can't log in
- `-p` to change the password for the account
- `-U` to unlock the account so that the user can log in

The `-L` parameter is especially handy. Use this to lock an account so that a user can't log in without having to remove the account and the user's data. To return the account to normal, just use the `-U` parameter.

passwd and chpasswd

A quick way to change just the password for a user is the `passwd` command:

```
# passwd test
Changing password for user test.
New UNIX password:
Retype new UNIX password:
passwd: all authentication tokens updated successfully.
#
```

If you just use the `passwd` command by itself, it will change your own password. Any user in the system can change his or her own password, but only the root user can change someone else's password.

The `-e` option is a handy way to force a user to change the password on the next log in. This allows you to set the user's password to a simple value and then force them to

change it to something harder that they can remember.

If you ever need to do a mass password change for lots of users on the system, the `chpasswd` command can be a lifesaver. The `chpasswd` command reads a list of login name and password pairs (separated by a colon) from the standard input automatically encrypts the password, and sets it for the user account. You can also use the redirection command to redirect a file of `userid:password` pairs into the command:

```
# chpasswd < users.txt
#
```

chsh, chfn, and chage

The `chsh`, `chfn`, and `chage` utilities are specialized for specific account modification functions. The `chsh` command allows you to quickly change the default login shell for a user. You must use the full pathname for the shell, and not just the shell name:

```
# chsh -s /bin/csh test
Changing shell for test.
Shell changed.
#
```

The `chfn` command provides a standard method for storing information in the comments field in the `/etc/passwd` file. Instead of just inserting random text, such as names or nicknames, or even just leaving the comment field blank, the `chfn` command uses specific information used in the Unix `finger` command to store information in the comment field. The `finger` command allows you to easily find information about people on your Linux system:

```
# finger rich
Login: rich                               Name: Rich Blum
Directory: /home/rich                     Shell: /bin/bash
On since Thu Sep 20 18:03 (EDT) on pts/0 from 192.168.1.2
No mail.
No Plan.
#
```

Note

Because of security concerns, many Linux system administrators disable the `finger` command on their systems.

If you use the `chfn` command with no parameters, it queries you for the appropriate values to enter in to the comment field:

```
# chfn test
Changing finger information for test.
Name []: Ima Test
Office []: Director of Technology
Office Phone []: (123)555-1234
Home Phone []: (123)555-9876
```

```
Finger information changed.
```

```
# finger test
Login: test                               Name: Ima Test
Directory: /home/test                     Shell: /bin/csh
Office: Director of Technology             Office Phone: (123)555-1234
Home Phone: (123)555-9876
Never logged in.
No mail.
No Plan.
#
```

If you now look at the entry in the `/etc/passwd` file, it looks like this:

```
# grep test /etc/passwd
test:x:504:504:Ima Test,Director of Technology,(123)555-
1234,(123)555-9876:/home/test:/bin/csh
#
```

All of the finger information is neatly stored away in the `/etc/passwd` file entry.

Finally, the `chage` command helps you manage the password aging process for user accounts. There are several parameters to set individual values, shown in [Table 6.4](#).

Table 6.4 The `chage` Command Parameters

Parameter	Description
<code>-d</code>	Set the number of days since the password was last changed.
<code>-E</code>	Set the date the password will expire.
<code>-I</code>	Set the number of days of inactivity after the password expires to lock the account.
<code>-m</code>	Set the minimum number of days between password changes.
<code>-W</code>	Set the number of days before the password expires that a warning message appears.

The `chage` date values can be expressed using one of two methods:

- A date in YYYY-MM-DD format
- A numerical value representing the number of days since January 1, 1970

One neat feature of the `chage` command is that it allows you to set an expiration date for an account. Using this feature, you can create temporary user accounts that automatically expire on a set date, without your having to remember to delete them! Expired accounts are similar to

locked accounts. The account still exists, but the user can't log in with it.

Using Linux Groups

User accounts are great for controlling security for individual users, but they aren't so good at allowing groups of users to share resources. To accomplish this, the Linux system uses another security concept, called *groups*.

Group permissions allow multiple users to share a common set of permissions for an object on the system, such as a file, directory, or device (more on that later in the “Decoding File Permissions” section).

Linux distributions differ somewhat on how they handle default group memberships. Some Linux distributions create just one group that contains all of the user accounts as members. You need to be careful if your Linux distribution does this, as your files may be readable by all other users on the system. Other distributions create a separate user account for each user to provide a little more security.

Each group has a unique GID, which, like UIDs, is a unique numerical value on the system. Along with the GID, each group has a unique group name. There are a few group utilities you can use to create and manage your own groups on the Linux system. This section discusses how group information is stored and how to use the group utilities to create new groups and modify existing groups.

The /etc/group File

Just like user accounts, group information is stored in a file on the system. The `/etc/group` file contains information about each group used on the system. Here are a few examples from a typical `/etc/group` file on a Linux system:

```
root:x:0:root
bin:x:1:root,bin,daemon
daemon:x:2:root,bin,daemon
sys:x:3:root,bin,adm
adm:x:4:root,adm,daemon
rich:x:500:
mama:x:501:
katie:x:502:
jessica:x:503:
mysql:x:27:
test:x:504:
```

Like UIDs, GIDs are assigned using a special format. Groups used for system accounts are assigned GIDs below 500, and user groups are assigned GIDs starting at 500. The `/etc/group` file uses four fields:

- The group name
- The group password
- The GID
- The list of user accounts that belong to the group

The group password allows a non-group member to temporarily become a member of the group by using the password. This feature is not used all that commonly, but it does exist.

You should never add users to groups by editing the `/etc/group` file. Instead, use the `usermod` command (discussed

earlier in the “Linux Security” section) to add a user account to a group. Before you can add users to different groups, you must create the groups.

Note

The list of user accounts is somewhat misleading. You'll notice that there are several groups in the list that don't have any users listed. This isn't because they don't have any members. When a user account uses a group as the default group in the `/etc/passwd` file, the user account doesn't appear in the `/etc/group` file as a member. This has caused confusion for more than one system administrator over the years!

Creating New Groups

The `groupadd` command allows you to create new groups on your system:

```
# /usr/sbin/groupadd shared
# tail /etc/group
haldaemon:x:68:
xfs:x:43:
gdm:x:42:
rich:x:500:
mama:x:501:
katie:x:502:
jessica:x:503:
mysql:x:27:
test:x:504:
shared:x:505:
#
```

When you create a new group, there are no users assigned to it by default. The `groupadd` command doesn't provide an option for adding user accounts to the group.

Instead, to add new users, use the `usermod` command:

```
# /usr/sbin/usermod -G shared rich
# /usr/sbin/usermod -G shared test
# tail /etc/group
haldaemon:x:68:
xfs:x:43:
gdm:x:42:
rich:x:500:
mama:x:501:
katie:x:502:
jessica:x:503:
mysql:x:27:
test:x:504:
shared:x:505:rich, test
#
```

The `shared` group now has two members, `test` and `rich`. The `-G` parameter in `usermod` appends the new group to the list of groups for the user account.

Note

If you change the user groups for an account that is currently logged into the system, the user will have to log out then back in for the group changes to take effect.

Caution

Be careful when assigning groups for user accounts. If you use the `-g` parameter, the group name you specify replaces the default group for the user account. The `-G` parameter adds the group to the list of groups the user belongs to, keeping the default group intact.

Modifying Groups

As you can see from the `/etc/group` file, there isn't too much information about a group for you to modify. The `groupmod` command allows you to change the GID (using the `-g` parameter) or the group name (using the `-n` parameter) of an existing group:

```
# /usr/sbin/groupmod -n sharing shared
# tail /etc/group
haldaemon:x:68:
xfs:x:43:
gdm:x:42:
rich:x:500:
mama:x:501:
katie:x:502:
jessica:x:503:
mysql:x:27:
test:x:504:
sharing:x:505:test,rich
#
```

When changing the name of a group, the GID and group members remain the same, only the group name changes. Because all security permissions are based on the GID, you can change the name of a group as often as you wish without adversely affecting file security.

Decoding File Permissions

Now that you know about users and groups, it's time to decode the cryptic file permissions you've seen when using the `ls` command. This section describes how to decipher the permissions and where they come from.

Using File Permission Symbols

If you remember from Chapter 3, the `ls` command allows you to see the file permissions for files, directories, and devices on the Linux system:

```
$ ls -l
total 68
-rw-rw-r-- 1 rich rich   50 2010-09-13 07:49 file1.gz
-rw-rw-r-- 1 rich rich   23 2010-09-13 07:50 file2
-rw-rw-r-- 1 rich rich   48 2010-09-13 07:56 file3
-rw-rw-r-- 1 rich rich   34 2010-09-13 08:59 file4
-rwxrwxr-x 1 rich rich 4882 2010-09-18 13:58 myprog
-rw-rw-r-- 1 rich rich  237 2010-09-18 13:58 myprog.c
drwxrwxr-x 2 rich rich 4096 2010-09-03 15:12 test1
drwxrwxr-x 2 rich rich 4096 2010-09-03 15:12 test2
$
```

The first field in the output listing is a code that describes the permissions for the files and directories. The first character in the field defines the type of the object:

- `-` for files
- `d` for directories
- `l` for links
- `c` for character devices
- `b` for block devices
- `n` for network devices

After that, there are three sets of three characters. Each set of three characters defines an access permission triplet:

- `r` for read permission for the object
- `w` for write permission for the object
- `x` for execute permission for the object

If a permission is denied, a dash appears in the location.

The three sets relate the three levels of security for the object:

- The owner of the object
- The group that owns the object
- Everyone else on the system

This is broken down in [Figure 6.1](#).

Figure 6.1 The Linux file permissions

```
-rwxrwxr-x 1 rich rich 4882 2010-09-18 13:58 myprog
```



permissions for everyone else

permissions for group members

permissions for the file owner

The easiest way to discuss this is to take an example and decode the file permissions one by one:

```
-rwxrwxr-x 1 rich rich 4882 2010-09-18 13:58 myprog
```

The file `myprog` has the following sets of permissions:

- `rwx` for the file owner (set to the login name `rich`)
- `rwx` for the file group owner (set to the group name `rich`)
- `r-x` for everyone else on the system

These permissions indicate that the user login name `rich` can read, write, and execute the file (considered full permissions). Likewise, members in the group `rich` can also read, write, and execute the file. However, anyone else

not in the rich group can only read and execute the file; the `w` is replaced with a dash, indicating that write permissions are not assigned to this security level.

Default File Permissions

You may be wondering about where these file permissions come from. The answer is *umask*. The `umask` command sets the default permissions for any file or directory you create:

```
$ touch newfile
$ ls -al newfile
-rw-r--r--  1 rich      rich          0 Sep 20 19:16 newfile
$
```

The `touch` command created the file using the default permissions assigned to my user account. The `umask` command shows and sets the default permissions:

```
$ umask
0022
$
```

Unfortunately, the `umask` command setting isn't overtly clear, and trying to understand exactly how it works makes things even muddier. The first digit represents a special security feature called the *sticky bit*. We'll talk more about that later on in this chapter in the "Sharing Files" section.

The next three digits represent the octal values of the `umask` for a file or directory. To understand how `umask` works, you first need to understand octal mode security settings.

Octal mode security settings take the three `rwX` permission values and convert them into a 3-bit binary value, represented by a single octal value. In the binary

representation, each position is a binary bit. Thus, if the read permission is the only permission set, the value becomes `r--`, relating to a binary value of 100, indicating the octal value of 4. [Table 6.5](#) shows the possible combinations you'll run into.

Table 6.5 Linux File Permission Codes

Permissions	Binary	Octal	Description
---	000	0	No permissions
--x	001	1	Execute-only permission
-w-	010	2	Write-only permission
-wx	011	3	Write and execute permissions
r--	100	4	Read-only permission
rx-	101	5	Read and execute permissions
rw-	110	6	Read and write permissions
rwX	111	7	Read, write, and execute permissions

Octal mode takes the octal permissions and lists three of them in order for the three security levels (user, group, and everyone). Thus, the octal mode value 664 represents read and write permissions for the user and group, but read-only permission for everyone else.

Now that you know about octal mode permissions, the `umask` value becomes even more confusing. The octal mode shown for the default `umask` on my Linux system is 0022, but the file I created had an octal mode permission of 644. How did that happen?

The `umask` value is just that, a mask. It masks out the permissions you don't want to give to the security level. Now we have to dive into some octal arithmetic to figure out

the rest of the story.

The umask value is subtracted from the full permission set for an object. The full permission for a file is mode 666 (read/write permission for all), but for a directory it's 777 (read/write/execute permission for all).

Thus, in the example, the file starts out with permissions 666, and the umask of 022 is applied, leaving a file permission of 644.

The umask value is normally set in the `/etc/profile` startup file (see Chapter 5). You can specify a different default umask setting using the `umask` command:

```
$ umask 026
$ touch newfile2
$ ls -l newfile2
-rw-r----- 1 rich      rich          0 Sep 20 19:46 newfile2
$
```

By setting the umask value to 026, the default file permissions become 640, so the new file now is restricted to read-only for the group members, and everyone else on the system has no permissions to the file.

The umask value also applies to making new directories:

```
$ mkdir newdir
$ ls -l
drwxr-x--x  2 rich      rich          4096 Sep 20 20:11 newdir/
$
```

Because the default permissions for a directory are 777, the resulting permissions from the umask are different from those of a new file. The 026 umask value is subtracted from 777, leaving the 751 directory permission setting.

Changing Security Settings

If you've already created a file or directory and need to change the security settings on it, there are a few different utilities available in Linux. This section shows you how to change the existing permissions, the default owner, and the default group settings for a file or directory.

Changing Permissions

The `chmod` command allows you to change the security settings for files and directories. The format of the `chmod` command is:

```
chmod options mode file
```

The `mode` parameter allows you to set the security settings using either octal or symbolic mode. The octal mode settings are pretty straightforward; just use the standard three-digit octal code you want the file to have:

```
$ chmod 760 newfile
$ ls -l newfile
-rwxrwx---    1 rich      rich              0 Sep 20 19:16 newfile
$
```

The octal file permissions are automatically applied to the file indicated. The symbolic mode permissions are not so easy to implement.

Instead of using the normal string of three sets of three characters, the `chmod` command takes a different approach. The following is the format for specifying a permission in symbolic mode:

```
[ugoa...][[+-=][rwxXstugo...]
```

Makes perfectly good sense, doesn't it? The first group of characters defines to whom the new permissions apply:

- `u` for the user
- `g` for the group
- `o` for others (everyone else)
- `a` for all of the above

Next, a symbol is used to indicate whether you want to add the permission to the existing permissions (+), subtract the permission from the existing permission (−), or set the permissions to the value (=).

Finally, the third symbol is the permission used for the setting. You may notice that there are more than the normal `rw` values here. The additional settings are:

- `x` to assign execute permissions only if the object is a directory or if it already had execute permissions
- `s` to set the UID or GID on execution
- `t` to save program text
- `u` to set the permissions to the owner's permissions
- `g` to set the permissions to the group's permissions
- `o` to set the permissions to the other's permissions

Using these permissions looks like this:

```
$ chmod o+r newfile
$ ls -l newfile
-rwxrw-r-- 1 rich rich 0 Sep 20 19:16 newfile
$
```

The `o+r` entry adds the read permission to whatever permissions the everyone security level already had.

```
$ chmod u-x newfile
$ ls -l newfile
```

The `u-x` entry removes the execute permission that the user already had. Note that the settings for the `ls` command indicate if a file has execution permissions by adding an asterisk to the file name.

The *options* parameters provide a few additional features to augment the behavior of the `chmod` command. The `-R` parameter performs the file and directory changes recursively. You can use wildcard characters for the file name specified, changing the permissions on multiple files with just one command.

Changing Ownership

Sometimes you need to change the owner of a file, such as when someone leaves an organization or a developer creates an application that needs to be owned by a system account when it's in production. Linux provides two commands for doing that. The `chown` command makes it easy to change the owner of a file, and the `chgrp` command allows you to change the default group of a file.

The format of the `chown` command is:

```
chown options owner[.group] file
```

You can specify either the login name or the numeric UID for the new owner of the file:

```
# chown dan newfile
# ls -l newfile
-rw-rw-r-- 1 dan rich 0 Sep 20 19:16 newfile
#
```

Simple. The `chown` command also allows you to change both the user and group of a file:

```
# chown dan.shared newfile
# ls -l newfile
-rw-rw-r--    1 dan      shared          0 Sep 20 19:16 newfile
#
```

If you really want to get tricky, you can just change the default group for a file:

```
# chown .rich newfile
# ls -l newfile
-rw-rw-r--    1 dan      rich            0 Sep 20 19:16 newfile
#
```

Finally, if your Linux system uses individual group names that match user login names, you can change both with just one entry:

```
# chown test. newfile
# ls -l newfile
-rw-rw-r--    1 test     test            0 Sep 20 19:16 newfile
#
```

The `chown` command uses a few different option parameters. The `-R` parameter allows you to make changes recursively through subdirectories and files, using a wildcard character. The `-h` parameter also changes the ownership of any files that are symbolically linked to the file.

Note

Only the root user can change the owner of a file. Any user can change the default group of a file, but the user must be a member of the groups the file is changed from and to.

The `chgrp` command provides an easy way to change just

the default group for a file or directory:

```
$ chgrp shared newfile
$ ls -l newfile
-rw-rw-r-- 1 rich shared 0 Sep 20 19:16 newfile
$
```

Now any member in the shared group can write to the file. This is one way to share files on a Linux system. However, sharing files among a group of people on the system can get tricky. The next section discusses how to do this.

Sharing Files

As you've probably already figured out, creating groups is the way to share access to files on the Linux system. However, for a complete file-sharing environment, things are more complicated.

As you've already seen in the “Decoding File Permissions” section, when you create a new file, Linux assigns the file permissions of the new file using your default UID and GID. To allow others access to the file, you need to either change the security permissions for the everyone security group or assign the file a different default group that contains other users.

This can be a pain in a large environment if you want to create and share documents among several people. Fortunately, there's a simple solution for how to solve this problem.

There are three additional bits of information that Linux stores for each file and directory:

- **The set user id (SUID):** When a file is executed by a user, the program runs under the permissions of the file owner.
- **The set group id (SGID):** For a file, the program runs under the permissions of the file group. For a directory, new files created in the directory use the directory group as the default group.
- **The sticky bit:** The file remains (sticks) in memory after the process ends.

The SGID bit is important for sharing files. By enabling the SGID bit, you can force all new files created in a shared directory to be owned by the directory's group and now the individual user's group.

The SGID is set using the `chmod` command. It's added to the beginning of the standard three-digit octal value (making a four-digit octal value), or you can use the symbol `s` in symbolic mode.

If you're using octal mode, you'll need to know the arrangement of the bits, shown in [Table 6.6](#).

Table 6.6 The chmod SUID, SGID, and Sticky Bit Octal Values

Binary	Octal	Description
000	0	All bits are cleared.
001	1	The sticky bit is set.
010	2	The SGID bit is set.

011	3	The SGID and sticky bits are set.
100	4	The SUID bit is set.
101	5	The SUID and sticky bits are set.
110	6	The SUID and SGID bits are set.
111	7	All bits are set.

So, to create a shared directory that always sets the directory group for all new files, all you need to do is set the SGID bit for the directory:

```
$ mkdir testdir
$ ls -l
drwxrwxr-x    2 rich      rich          4096 Sep 20 23:12 testdir/
$ chgrp shared testdir
$ chmod g+s testdir
$ ls -l
drwxrwsr-x    2 rich      shared        4096 Sep 20 23:12 testdir/
$ umask 002
$ cd testdir
$ touch testfile
$ ls -l
total 0
-rw-rw-r--    1 rich      shared          0 Sep 20 23:13 testfile
$
```

The first step is to create a directory that you want to share using the `mkdir` command. Next, the `chgrp` command is used to change the default group for the directory to a group that contains the members who need to share files. Finally, the SGID bit is set for the directory to ensure that any files created in the directory use the shared group name as the default group.

For this environment to work properly, all of the group members need to have their umask values set to make files writable by group members. In the preceding example, the `umask` is changed to `002` so that the files are writable by the

group.

After all that's done, any member of the group can go to the shared directory and create a new file. As expected, the new file uses the default group of the directory, not the user account's default group. Now any user in the shared group can access this file.

Summary

This chapter discussed the command line commands you need to know to manage the Linux security on your system. Linux uses a system of user IDs and group IDs to protect access to files, directories, and devices. Linux stores information about user accounts in the `/etc/passwd` file and information about groups in the `/etc/group` file. Each user is assigned a unique numeric user ID, along with a text login name to identify the user in the system. Groups are also assigned unique numerical group IDs, and text group names. A group can contain one or more users to allowed shared access to system resources.

There are several commands available for managing user accounts and groups. The `useradd` command allows you to create new user accounts, and the `groupadd` command allows you to create new group accounts. To modify an existing user account, use the `usermod` command. Similarly, the `groupmod` command is used to modify group account information.

Linux uses a complicated system of bits to determine

access permissions for files and directories. Each file contains three security levels of protection: the file's owner, a default group that has access to the file, and a level for everyone else on the system. Each security level is defined by three access bits: read, write, and execute. The combination of three bits is often referred to by the symbols `rwX`, for read, write, and execute. If a permission is denied, its symbol is replaced with a dash (such as `r--` for read-only permission).

The symbolic permissions are often referred to as octal values, with the three bits combined into one octal value and three octal values representing the three security levels. The `umask` command is used to set the default security settings for files and directories created on the system. The system administrator normally sets a default `umask` value in the `/etc/profile` file, but you can use the `umask` command to change your `umask` value at any time.

The `chmod` command is used to change security settings for files and directories. Only the file's owner can change permissions for a file or directory. However, the root user can change the security settings for any file or directory on the system. The `chown` and `chgrp` commands can be used to change the default owner and group of the file.

Finally, the chapter closed out with a discussion on how to use the set GID bit to create a shared directory. The SGID bit forces any new files or directories created in a directory to use the default group name of the parent directory, not that of the user who created them. This

provides an easy way to share files between users on the system.

Now that you're up to speed with file permissions, it's time to take a closer look at how to work with the actual filesystem in Linux. The next chapter shows you how to create new partitions in Linux from the command line and then how to format the new partitions so that they can be used in the Linux virtual directory.

Chapter 7

Managing Filesystems

In This Chapter

- What is a filesystem?
- Linux filesystems
- Filesystem commands

When you're working with your Linux system, one of the decisions you'll need to make is what filesystems to use for the storage devices. Most Linux distributions kindly provide a default filesystem for you at installation time, and most beginning Linux users just use that by default without giving the topic another thought.

While using the default filesystem choice isn't necessarily a bad thing, sometimes it helps to know the other options available to you. This chapter discusses the different filesystem options you have available in the Linux world and shows you how to create and manage them from the Linux command line.

Exploring Linux Filesystems

Chapter 3 discussed how Linux uses a *filesystem* to store files and folders on a storage device. The filesystem provides a way for Linux to bridge the gap between the ones and zeroes stored in the hard drive and the files and folders you work with in your applications.

Linux supports several types of filesystems to manage files and folders. Each filesystem implements the virtual directory structure on storage devices using slightly different features. This section walks you through the strengths and weaknesses of the more common filesystems used in the Linux environment.

The Basic Linux Filesystems

The original Linux system used a simple filesystem that mimicked the functionality of the Unix filesystem. This section discusses the evolution of that filesystem.

The ext Filesystem

The original filesystem introduced with the Linux operating system is called the *extended filesystem* (or just ext for short). It provides a basic Unix-like filesystem for Linux, using virtual directories to handle physical devices, and storing data in fixed-length blocks on the physical devices.

The ext filesystem uses a system called *inodes* to track information about the files stored in the virtual directory. The

inode system creates a separate table on each physical device, called the *inode table*, to store the file information. Each stored file in the virtual directory has an entry in the inode table. The *extended* part of the name comes from the additional data that it tracks on each file, which consists of:

- The file name
- The file size
- The owner of the file
- The group the file belongs to
- Access permissions for the file
- Pointers to each disk block that contains data from the file

Linux references each inode in the inode table using a unique number (called the *inode number*), assigned by the filesystem as data files are created. The filesystem uses the inode number to identify the file rather than having to use the full file name and path.

The ext2 Filesystem

The original ext filesystem had quite a few limitations, such as limiting files to only 2GB in size. Not too long after Linux was first introduced, the ext filesystem was upgraded to create the second extended filesystem, called *ext2*.

As you can guess, the ext2 filesystem is an expansion of the basic abilities of the ext filesystem, but maintains the same structure. The ext2 filesystem expands the inode table format to track additional information about each file

on the system.

The ext2 inode table adds the created, modified, and last accessed time values for files to help system administrators track file access on the system. The ext2 filesystem also increases the maximum file size allowed to 2TB (then in later versions of ext2, that was increased to 32TB) to help accommodate large files commonly found in database servers.

In addition to expanding the inode table, the ext2 filesystem also changed the way in which files are stored in the data blocks. A common problem with the ext filesystem was that as a file is written to the physical device, the blocks used to store the data tend to be scattered throughout the device (called *fragmentation*). Fragmentation of data blocks can reduce the performance of the filesystem, as it takes longer to search the storage device to access all of the blocks for a specific file.

The ext2 filesystem helps reduce fragmentation by allocating disk blocks in groups when you save a file. By grouping the data blocks for a file, the filesystem doesn't have to search all over the physical device for the data blocks to read the file.

The ext2 filesystem was the default filesystem used in Linux distributions for many years, but it, too, had its limitations. The inode table, while a nice feature that allows the filesystem to track additional information about files, can cause problems that can be fatal to the system. Each time the filesystem stores or updates a file, it has to modify

the inode table with the new information. The problem is that this isn't always a fluid action.

If something should happen to the computer system between the file being stored and the inode table being updated, the two would become out of sync. The ext2 filesystem is notorious for easily becoming corrupted due to system crashes and power outages. Even if the file data is stored just fine on the physical device, if the inode table entry wasn't completed, the ext2 filesystem wouldn't even know that the file existed!

It wasn't long before developers were exploring a different avenue of Linux filesystems.

Journaling Filesystems

Journaling filesystems provide a new level of safety to the Linux system. Instead of writing data directly to the storage device and then updating the inode table, journaling filesystems write file changes into a temporary file (called the *journal*) first. After data is successfully written to the storage device and the inode table, the journal entry is deleted.

If the system should crash or suffer a power outage before the data can be written to the storage device, the journaling filesystem just reads through the journal file and processes any uncommitted data left over.

There are three different methods of journaling commonly used in Linux, each with different levels of protection. These

are shown in [Table 7.1](#).

Table 7.1 Journaling Filesystem Methods

Method	Description
Data mode	Both inode and file data are journaled. Low risk of losing data, but poor performance.
Ordered mode	Only inode data written to the journal, but not removed until file data is successfully written. Good compromise between performance and safety.
Writeback mode	Only inode data written to the journal, no control over when the file data is written. Higher risk of losing data, but still better than not using journaling.

The data mode journaling method is by far the safest for protecting data, but it is also the slowest. All of the data written to a storage device must be written twice, once to the journal, then again to the actual storage device. This can cause poor performance, especially for systems that do a lot of data writing.

Over the years, a few different journaling filesystems have appeared in Linux. The following sections describe the popular Linux journaling filesystems available.

The Extended Linux Journaling Filesystems

The same group that developed the ext and ext2 filesystems as part of the Linux project also created journaling versions of the filesystems. These journaling filesystems are compatible with the ext2 filesystem, and it's

easy to convert back and forth between them. There are currently two separate journaling filesystems based on the ext2 filesystem.

The ext3 Filesystem

The *ext3* filesystem was added to the Linux kernel in 2001, and up until recently was the default filesystem used by just about all Linux distributions. It uses the same inode table structure as the *ext2* filesystem, but adds a journal file to each storage device to journal the data written to the storage device.

By default, the *ext3* filesystem uses the ordered mode method of journaling, only writing the inode information to the journal file, but not removing it until the data blocks have been successfully written to the storage device. You can change the journaling method used in the *ext3* filesystem to either data or writeback modes with a simple command line option when creating the filesystem.

While the *ext3* filesystem added basic journaling to the Linux filesystem, there were still a few things it lacked. For example, the *ext3* filesystem doesn't provide any recovery from accidental deletion of files, there's no built-in data compression available (although there is a patch that can be installed separately that provides this feature), and the *ext3* filesystem doesn't support encrypting files. For those reasons developers in the Linux project chose to continue work on improving the *ext3* filesystem.

The ext4 Filesystem

The result of expanding the ext3 filesystem was (as you probably guessed) the *ext4* filesystem. The ext4 filesystem was officially supported in the Linux kernel in 2008, and is now the default filesystem used in most popular Linux distributions, such as Fedora and Ubuntu.

In addition to supporting compression and encryption, the ext4 filesystem also supports a feature called *extents*. Extents allocate space on a storage device in blocks, and only store the starting block location in the inode table. This helps save space in the inode table by not having to list all of the data blocks used to store data from the file.

The ext4 filesystem also incorporates *block preallocation*. If you want to reserve space on a storage device for a file that you know will grow in size, with the ext4 filesystem it's possible to allocate all of the expected blocks for the file, not just the blocks that physically exist. The ext4 filesystem fills in the reserved data blocks with zeroes, and knows not to allocate them for any other file.

The Reiser Filesystem

In 2001, Hans Reiser created the first journaling filesystem for Linux, called *ReiserFS*. The ReiserFS filesystem only supports writeback journaling mode, writing only the inode table data to the journal file. Because it writes only the inode table data to the journal, the ReiserFS filesystem is one of the fastest journaling filesystems in Linux.

Two interesting features incorporated into the ReiserFS filesystem are that you can resize an existing filesystem while it's still active, and that it uses a technique called *tailpacking*, which stuffs data from one file into empty space in a data block from another file. The active filesystem resizing feature is great if you have to expand an already created filesystem to accommodate more data.

Warning

Because Hans Reiser's legal issues, the status of the ReiserFS filesystem is in question. While the ReiserFS project is an open source project, Hans was the lead developer for the project and is now in jail. No recent development has been done on the ReiserFS filesystem, and it's not clear if anyone else will pick up development work for the project. Given that other alternatives for journaling filesystems are now available, most new Linux installations don't use the ReiserFS filesystem. However, you may still run into an existing Linux system that uses it.

The Journaled Filesystem

Possibly one of the oldest journaling filesystems around, the *Journaled File System* (JFS) was developed by IBM in 1990 for its AIX flavor of Unix. However, it wasn't until its second version that it was ported to the Linux environment.

Note

The official IBM name of the second version of the

JFS filesystem is JFS2, but most Linux systems refer to it as just JFS.

The JFS filesystem uses the ordered journaling method, storing only the inode table data in the journal, but not removing it until the actual file data is written to the storage device. This method is a compromise between the speed of the ReiserFS and the integrity of the data mode journaling method.

The JFS filesystem uses extent-based file allocation, allocating a group of blocks for each file written to the storage device. This method provides for less fragmentation on the storage device.

Outside of the IBM Linux offerings, the JFS filesystem isn't popularly used, but you may run into it in your Linux journey.

The XFS Filesystem

The XFS journaling filesystem is yet another filesystem originally created for a commercial Unix system that made its way into the Linux world. Silicon Graphics Incorporated (SGI) originally created XFS in 1994 for its commercial IRIX Unix system. It was released to the Linux environment for common use in 2002.

The XFS filesystem uses the writeback mode of journaling, which provides high performance but does introduce an amount of risk because the actual data isn't stored in the journal file. The XFS filesystem also allows

online resizing of the filesystem, similar to the ReiserFS filesystem, except XFS filesystems can only be expanded and not shrunk.

Working with Filesystems

Linux provides a few different utilities that make it easier to work with filesystems from the command line. You can add new filesystems or change existing filesystems from the comfort of your own keyboard. This section walks you through the commands for interacting with filesystems from a command line environment.

Creating Partitions

To start out, you'll need to create a *partition* on the storage device to contain the filesystem. The partition can be an entire disk, or it can be a subset of a disk, that contains a portion of the virtual directory.

The `fdisk` utility is used to help you organize partitions on any storage device installed on the system. The `fdisk` command is an interactive program that allows you to enter commands to walk through the steps of partitioning a hard drive.

To start the `fdisk` command, you'll need to specify the device name of the storage device you want to partition:

```
$ sudo fdisk /dev/sdc
[sudo] password for rich:
Device contains neither a valid DOS partition table, nor Sun, SGI or
```

```
OSF disklabel
Building a new DOS disklabel with disk identifier 0x4beedc66.
Changes will remain in memory only, until you decide to write them.
After that, of course, the previous content won't be recoverable.

Warning: invalid flag 0x0000 of partition table 4 will be corrected
by w(rite)

WARNING: DOS-compatible mode is deprecated. It's strongly recommended
        To switch off the mode (command 'c') and change display
        units to sectors (command 'u').

Command (m for help):
```

If this is the first time you're partitioning the storage device, fdisk will give you a warning that there isn't a partition table on the device.

The fdisk interactive command prompt uses single letter commands to instruct fdisk what to do. [Table 7.2](#) shows the commands available at the fdisk command prompt.

[Table 7.2](#) The fdisk Commands

Command	Description
a	Toggles a flag indicating if the partition is bootable
b	Edits the disklabel used by BSD Unix systems
c	Toggles the DOS compatibility flag
d	Deletes the partition
l	Lists the available partition types
m	Displays the command options
n	Adds a new partition
o	Creates a DOS partition table
p	Displays the current partition table
q	Quits without saving changes
s	Creates a new disklabel for Sun Unix systems

t	Changes the partition system ID
u	Changes the storage units used
v	Verifies the partition table
w	Writes the partition table to the disk
x	Advanced functions

While this list may look intimidating, there are usually just a few basic commands that you'll normally use in day-to-day work.

For starters, you can display the details of a storage device using the `p` command:

```
Command (m for help): p
```

```
Disk /dev/sdc: 5368 MB, 5368946688 bytes
255 heads, 63 sectors/track, 652 cylinders
Units = cylinders of 16065 * 512 = 8225280 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk identifier: 0x4beedc66
```

Device	Boot	Start	End	Blocks	Id	System
--------	------	-------	-----	--------	----	--------

```
Command (m for help):
```

The output shows that the storage device has 5368MB of space on it (5GB). The listing under the storage device details shows if there are any existing partitions on the device. The listing in this example doesn't show any partitions, so the device is not partitioned yet.

Next, you'll want to create a new partition on the storage device. Use the `n` command for that:

```
Command (m for help): n
```

```
Command action
```

```
e    extended
```

```
p    primary partition (1-4)
```

```
p
```

```
Partition number (1-): 1
First cylinder (1-652, default 1): 1
Last cylinder, +cylinders or +size{K,M,G} (1-652, default 652): +2G

Command (m for help):
```

Partitions can be created as either a *primary partition* or an *extended partition*. Primary partitions can be formatted with a filesystem directly, whereas extended partitions can only contain other primary partitions. The reason for extended partitions is that there can only be four partitions on a single storage device. You can extend that by creating multiple extended partitions, then creating primary partitions inside the extended partitions. This example creates a primary storage device, assigns it as partition number 1 on the storage device, and then allocates 2GB of the storage device space to it. You can see the results using the `p` command again:

```
Command (m for help): p

Disk /dev/sdc: 5368 MB, 5368946688 bytes
255 heads, 63 sectors/track, 652 cylinders
Units = cylinders of 16065 * 512 = 8225280 bytes
Sector size (logical/physical): 512 bytes / 512 bytes
I/O size (minimum/optimal): 512 bytes / 512 bytes
Disk identifier: 0x4beedc66
```

Device	Boot	Start	End	Blocks	Id	System
/dev/sdc1		1	262	2104483+	83	Linux

```
Command (m for help):
```

Now in the output there's a partition on the storage device (called `/dev/sdc1`). The `Id` entry defines how Linux treats the partition. There are lots of partition types that `fdisk` allows you to create. Using the `l` command lists the different types available. The default is type 83, which

defines a Linux filesystem. If you want to create a partition for a different filesystem (such as a Windows NTFS partition), just select a different partition type.

You can repeat the process to allocate the remaining space on the storage device to another Linux partition. After you've created the partitions you want, use the `w` command to save the changes to the storage device:

```
Command (m for help): w
The partition table has been altered!

Calling ioctl() to re-read partition table.
Syncing disks.
$
```

Now that you have a partition on the storage device, you're ready to format it with a Linux filesystem.

Note

Sometimes the hardest part of creating a new disk partition is trying to find the physical disk on your Linux system. Linux uses a standard format for assigning device names to hard drives, but you need to be familiar with the format. For older IDE drives, Linux uses `/dev/hdx`, where `x` is a letter based on the order the drive is detected (a for the first drive, b for the second, and so on). For both the newer SATA drives, and SCSI drives, Linux uses `/dev/sdx`, where `x` is a letter based on the order the drive is detected (again, a for the first drive, b for the second, and so on). It's always a good idea to double-check to make sure you are referencing the correct drive before formatting the partition!

Creating a Filesystem

Before you can store data on the partition, you must format it with a filesystem so that Linux can use it. Each filesystem type uses its own command line program to format partitions. [Table 7.3](#) lists the utilities used for the different filesystems discussed in this chapter.

Table 7.3 Command Line Programs to Create Filesystems

Utility	Purpose
<code>mkefs</code>	Create an ext filesystem.
<code>mke2fs</code>	Create an ext2 filesystem.
<code>mkfs.ext3</code>	Create an ext3 filesystem.
<code>mkfs.ext4</code>	Create an ext4 filesystem.
<code>mkreiserfs</code>	Create a ReiserFS filesystem.
<code>jfs_mkfs</code>	Create a JFS filesystem.
<code>mkfs.xfs</code>	Create an XFS filesystem.

Each filesystem command has lots of command line options that allow you to customize just how the filesystem is created in the partition. To see all of the command line options available, use the `man` command to display the manual pages for the filesystem command (see Chapter 2). All of the filesystem commands allow you to create a default filesystem with just the simple command with no options:

```
$ sudo mkfs.ext4 /dev/sdcl
[sudo] password for rich:
mke2fs 1.41.11 (14-Mar-2010)
Filesystem label=
OS type: Linux
Block size=4096 (log=2)
Fragment size=4096 (log=2)
```

```
Stride=0 blocks, Stripe width=0 blocks
131648 inodes, 526120 blocks
26306 blocks (5.00%) reserved for the super user
First data block=0
Maximum filesystem blocks=541065216
17 block groups
32768 blocks per group, 32768 fragments per group
7744 inodes per group
Superblock backups stored on blocks:
    32768, 98304, 163840, 229376, 294912

Writing inode tables: done
Creating journal (16384 blocks): done
Writing superblocks and filesystem accounting information: done

This filesystem will be automatically checked every 34 mounts or
180 days, whichever comes first.  Use tune2fs -c or -i to override.
$
```

The new filesystem uses the ext4 filesystem type, which is the journaling filesystem in Linux. Notice that part of the creation process was to create the new journal.

After you create the filesystem for the partition, the next step is to mount it on a mount point in the virtual directory so you can store data in the new filesystem. You can mount the new filesystem anywhere in your virtual directory where you need the extra space.

```
$ sudo mkdir /mnt/testing
$ sudo mount -t ext4 /dev/sdb1 /mnt/testing
$ ls -al /mnt/testing
total 24
drwxr-xr-x 3 root root 4096 2010-09-25 19:25 .
drwxr-xr-x 3 root root 4096 2010-09-25 19:38 ..
drwx----- 2 root root 16384 2010-09-25 19:25 lost+found
$
```

The `mkdir` command creates the mount point in the virtual directory, and the `mount` command adds the new hard drive partition to the mount point. Now you can save new files and folders on the new partition!

Warning

This method of mounting a filesystem only temporarily mounts the filesystem. When you reboot your Linux system, the filesystem won't be automatically mounted. To force Linux to automatically mount the new filesystem at boot time, add the new filesystem to the `/etc/fstab` file.

This chapter showed how to handle filesystems contained in physical storage devices. Linux also provides a couple of different ways to create logical storage devices for filesystems. The next section examines how you can use a logical storage device for your filesystems.

If Things Go Wrong

Even with fancy journaling filesystems, things can go wrong if power is unexpectedly lost, or if a wayward application locks up the system while file access is in progress. Fortunately, there are some command line tools available to help you make an attempt to restore the filesystem back to order.

Each filesystem has its own recovery command for interacting with the filesystem. That has the potential of getting ugly, as more and more filesystems are available in the Linux environment, making for lots of individual commands you have to know. Fortunately, there's a common front-end program available that can determine the filesystem on the storage device and use the

appropriate filesystem recovery command based on the filesystem being recovered.

The `fsck` command is used to check and repair any type of Linux filesystem, including all of the ones discussed earlier in this chapter — `ext`, `ext2`, `ext3`, `ext4`, `ReiserFS`, `JFS`, and `XFS`. The format of the command is:

```
fsck options filesystem
```

You can list multiple *filesystem* entries on the command line to check. Filesystems can be referenced using either the device name, the mount point in the virtual directory, or a special Linux UUID value assigned to the filesystem.

The `fsck` command uses the `/etc/fstab` file to automatically determine the filesystem on a storage device that's normally mounted on the system. If the storage device isn't normally mounted (such as if you just created a filesystem on a new storage device, you'll need to use the `-t` command line option to specify the filesystem type. [Table 7.4](#) lists the other command line options available.

Table 7.4 The `fsck` Command Line Options

Option	Description
-a	Automatically repair the filesystem if errors are detected.
-A	Check all of the filesystems listed in the <code>/etc/fstab</code> file.
-C	Display a progress bar for filesystems that support that feature (only <code>ext2</code> and <code>ext3</code>).
-N	Don't run the check, only display what checks would be performed.
-r	Prompt to fix if errors found.

-R	Skip the root filesystem if using the -A option.
-S	If checking multiple filesystems, perform the checks one at a time.
-t	Specify the filesystem type to check.
-T	Don't show the header information when starting.
-v	Produce verbose output during the checks.
-y	Automatically repair the filesystem if errors detected.

You may notice that some of the command line options are redundant. That's part of the problem of trying to implement a common front-end for multiple commands. Some of the individual filesystem repair commands have additional options that can be used. If you need to do more advanced error checking, you'll need to check the man pages for the individual filesystem repair tool to see if there are extended options specific to that filesystem.

Warning

You can run the `fsck` command on unmounted filesystems only. For most filesystems, you can just unmount the filesystem to check it, then remount it when done. However, because the root filesystem contains all of the core Linux commands and log files, you can't unmount it on a running system.

This is a time where having a Linux LiveCD comes in handy! Just boot your system with the LiveCD, and then run the `fsck` command on the root filesystem!

Logical Volume Managers

If you create your filesystems using standard partitions on hard drives, trying to add additional space to an existing filesystem can be somewhat of a painful experience. You can only expand a partition to the extent of the available space on the same physical hard drive. If there isn't any more space available on that hard drive, you're stuck having to get a larger hard drive and manually moving the existing filesystem to the new drive.

What would come in handy is a way to dynamically add more space to an existing filesystem by just adding a partition from another hard drive to the existing filesystem. The *Linux Logical Volume Manager* (LVM) software package allows you to do just that. It provides an easy way for you to manipulate disk space on a Linux system without having to rebuild entire filesystems.

The Logical Volume Management Layout

The core of logical volume management is how it handles the physical hard drive partitions installed on the system. In the logical volume management world, hard drives are called *physical volumes* (PV). Each PV maps to a specific physical partition created on a hard drive.

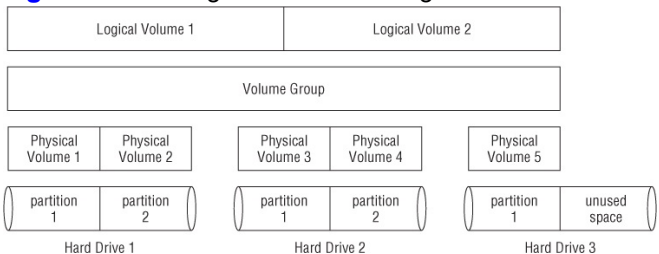
Multiple PV elements are pooled together to create a *volume group* (VG). The logical volume management

system treats the VG like a physical hard drive, but in reality the VG may consist of multiple physical partitions spread across multiple hard drives. The VG provides a platform to create the logical partitions, which actually contain the filesystem.

The final layer in the structure is the *logical volume (LV)*. The LV creates the partition environment for Linux to create a filesystem, acting similar to a physical hard disk partition as far as Linux is concerned. The Linux system treats the LV just like a physical partition. You can format the LV using any one of the standard Linux filesystems, and then add it to the Linux virtual directory at a mount point.

[Figure 7.1](#) shows the basic layout of a typical Linux logical volume management environment.

[Figure 7.1](#) The Logical Volume Management environment



The volume group shown in [Figure 7.1](#) spans across three separate physical hard drives, covering five separate physical partitions. Inside the volume group are two separate logical volumes. The Linux system treats each

logical volume just like a physical partition. Each logical volume can be formatted as an ext4 filesystem, and then mounted to a specific location in the virtual directory.

Notice in [Figure 7.1](#) that the third physical hard drive has an unused partition. Using logical volume management, you can easily assign this unused partition to the existing volume group at a later time, and then either use it to create a new logical volume or add it to expand one of the existing logical volumes when you need more space.

Likewise, if you add a new hard drive to the system, the local volume management system allows you to add it to the existing volume group, and then create more space for one of the existing logical volumes, or start a new logical volume to be mounted. That's a much better way of handling expanding filesystems!

The LVM in Linux

The Linux LVM was developed by Heinz Mauelshagen and released to the Linux community in 1998. It allows you to manage a complete logical volume management environment in Linux using simple command line commands.

Two versions of Linux LVM are available:

- **LVM1:** The original LVM package released in 1998, and available in only the 2.4 Linux kernels. It provides only basic Logical Volume Management features.
- **LVM2:** An updated version of the LVM, available in

the 2.6 Linux kernels. It provides additional features over the standard LVM1 features.

Most modern Linux distributions using the 2.6 kernel version provide support for LVM2. Besides the standard logical volume management features, LVM2 provides a few other nice things for you to use in your Linux system.

Snapshots

The original Linux LVM allows you to copy an existing logical volume to another device while the logical volume is active. This feature is called a *snapshot*. Snapshots are great for backing up important data that can't be locked due to high availability requirements. Traditional backup methods usually lock files as they're being copied to the backup media. The snapshot allows you to continue running mission critical Web or database servers while performing the copy. Unfortunately, LVM1 allows you to create only a read-only snapshot. Once you create the snapshot you can't write to it.

LVM2 allows you to create a read-write snapshot of an active logical volume. With the read-write copy, you can remove the original logical volume and mount the snapshot as a replacement. This feature is great for fast fail-overs, or for experimenting with applications that modify data that may need to be restored if something fails.

Striping

Another interesting feature that LVM2 provides is *striping*. With striping, a logical volume is created across multiple physical hard drives. When the Linux LVM writes a file to the logical volume, the data blocks in the file are spread across the multiple hard drives. Each successive block of data is written to the next hard drive.

Striping helps improve disk performance, as Linux can write the multiple data blocks for a file to the multiple hard drives simultaneously, rather than having to wait for a single hard drive to move the read/write head to different locations. This improvement also applies to reading sequentially accessed files, as the LVM can read data from the multiple hard drives simultaneously.

Warning

LVM striping is not the same as RAID striping. LVM striping doesn't provide a parity entry, which creates the fault-tolerant environment. In fact, LVM striping may increase the chance of a file being lost due to a hard drive failure. A single disk failure can result in multiple logical volumes being inaccessible.

Mirroring

Just because you install a filesystem using LVM doesn't mean that things can't still go wrong in the filesystem. Just as in a physical partition, LVM logical volumes are susceptible to power outages and disk crashes. Once a

filesystem becomes corrupt, there's always a possibility that you won't be able to recover it.

The LVM snapshot process provides some comfort knowing that you can create a backup copy of a logical volume at any time, but for some environments that may not be enough. Systems that have lots of data changes, such as database servers, may store hundreds or thousands of records since the last snapshot.

A solution to this problem is the LVM *mirror*. A mirror is a complete copy of a logical volume that's updated in real time. When you create the mirror logical volume, LVM synchronizes the original logical volume to the mirror copy. Depending on the size of the original logical volume, this may take some time to complete.

Once the original synchronization is complete, LVM performs two writes for each write process in the filesystem — one to the main logical volume and one to the mirrored copy. As you can guess, this process does slow down write performance on the system. However, if the original logical volume should become corrupt for some reason, you have a complete up-to-date copy at your fingertips!

Using the Linux LVM

Now that you've seen what the Linux LVM can do, this section discusses how to implement it to help organize the disk space on your system. The Linux LVM package only provides command line programs for creating and managing all of the components in the Logical Volume

Management system. Some Linux distributions include graphical front-ends to the command line commands, but for complete control of your LVM environment, it's best to get comfortable working directly with the commands.

Defining Physical Volumes

The first step in the process is to convert the physical partitions on the hard drive into physical volume extents used by the Linux LVM. Our friend the `fdisk` command will help us here. After creating the basic Linux partition, you need to change the partition type using the `t` command:

```
Command (m for help): t
Selected partition 1
Hex code (type L to list codes): 8e
Changed system type of partition 1 to 8e (Linux LVM)
```

```
Command (m for help): p
```

```
Disk /dev/sdc: 2147 MB, 2147992064 bytes
255 heads, 63 sectors/track, 261 cylinders
Units = cylinders of 16065 * 512 = 8225280 bytes
Disk identifier: 0x4bc26133
```

Device	Boot	Start	End	Blocks	Id	System
/dev/sdc1		1	261	2096451	8e	Linux LVM

```
Command (m for help):
```

The `8e` partition type denotes that the partition will be used as part of a Linux LVM system and not as a direct filesystem, as you saw with the `83` partition type earlier.

The next step is to use the partition to create the actual physical volume. That's done using the `pvcreate` command:

```
$ sudo pvcreate /dev/sdc1
Physical volume "/dev/sdc1" successfully created
$
```

The `pvcreate` command defines the physical volume to use for the PV. It simply tags the partition as a physical volume in the Linux LVM system. You can use the `pvdisplay` command to display a list of physical volumes you've created if you'd like to see your progress along the way.

Creating Volume Groups

The next step in the process is to create one or more volume groups from the physical volumes. There are no set rules for how many volume groups you need to create for your system — you can add all of the available physical volumes to a single volume group, or you can create multiple volume groups by combining different physical volumes.

To create the volume group from the command line, you need to use the `vgcreate` command. The `vgcreate` command requires a few command line parameters to define the volume group name, as well as the name of the physical volumes you're using to create the volume group:

```
$ sudo vgcreate Voll /dev/sdc1
Volume group "Voll" successfully created
$
```

That's not all too exciting for output! If you'd like to see some details about the newly created volume group, use the `vgdisplay` command:

```
$ sudo vgdisplay
--- Volume group ---
VG Name                Voll
System ID
Format                 lvm2
Metadata Areas         1
```


Metadata Sequence No	1
VG Access	read/write
VG Status	resizable
MAX LV	0
Cur LV	0
Open LV	0
Max PV	0
Cur PV	1
Act PV	1
VG Size	2.00 GB
PE Size	4.00 MB
Total PE	511
Alloc PE / Size	0 / 0
Free PE / Size	511 / 2.00 GB
VG UUID	Cy1lHZ-Y840-8TUn-Wvti-4S6Q-bHHT-C113I0

\$

This example creates a volume group named `Vol1`, using the physical volume created on the `/dev/sdc1` partition.

Now that you have one or more volume groups created, you're ready to create the logical volume.

Creating Logical Volumes

The logical volume is what the Linux system uses to emulate a physical partition, and holds the filesystem. The Linux system handles the logical volumes just like a physical partition, allowing you to define filesystems in the logical volume and then mount the filesystem into the virtual directory.

To create the logical volume, use the `lvcreate` command. While you can usually get away without using command line options in the other Linux LVM commands, the `lvcreate` command requires at least some options to be entered. [Table 7.5](#) shows the available command line options.

Table 7.5 The lvcreate Options

Option	Long Option Name	Description
-c	-- chunksize	Specify the chunksize of the snapshot logical volume.
-C	-- contiguous	Set or reset the contiguous allocation policy.
-i	--stripes	Specify the number of stripes.
-I	-- stripesize	Specify the size of each stripe.
-l	--extents	Specify the number of logical extents to allocate to a new logical volume, or the percent of the logical extents to use.
-L	--size	Specify the disk size to allocate to a new logical volume.
	--minor	Specify the minor number of the device.
-m	--mirrors	Create a mirrored logical volume.
-M	-- persistent	Make the minor number persistent.
-n	--name	Specify the name of the new logical volume.
-p	-- permission	Set read/write permission for the logical volume.
-r	-- readahead	Set read ahead sector count.
-R	-- regionsize	Specify the size to divide the mirror regions into.
-s	--snapshot	Create a snapshot logical volume.
-Z	--zero	Set the first 1KB of data on the new logical volume to zeros.

While the command line options may look intimidating, for most situations you can get by with a minimal amount of options:

```
$ sudo lvcreate -l 100%FREE -n lvtest Vol1
Logical volume "lvtest" created
$
```

If you want to see the details of what you created, use the `lvdisplay` command:

```
$ sudo lvdisplay
--- Logical volume ---
LV Name                /dev/Vol1/lvtest
VG Name                Vol1
LV UUID                usDxti-pAEj-fEIz-3kWV-LNAu-PFNx-2LGqNv
LV Write Access        read/write
LV Status              available
# open                 0
LV Size                2.00 GB
Current LE             511
Segments              1
Allocation             inherit
Read ahead sectors     auto
- currently set to    256
Block device          253:2
$
```

Now you can see just what you created! Notice that the volume group name (`Vol1`) is used to identify the volume group to use when creating the new logical volume.

The `-l` parameter defines how much of the available space on the volume group specified to use for the logical volume. Notice that you can specify the value as a percent of the free space in the volume group. This example used all (100%) of the free space for the new logical volume.

You can use the `-l` parameter to specify the size as a percentage of the available space, or the `-L` parameter to specify the actual size in bytes, kilobytes (KB), megabytes

(MB), or gigabytes (GB). The `-n` parameter allows you to provide a name for the logical volume (called `lvtest` in this example).

Creating the Filesystem

After you run the `lvcreate` command, the logical volume exists but doesn't have a filesystem. To do that, you'll need to use the appropriate command line program for the filesystem you want to create:

```
$ sudo mkfs.ext4 /dev/Vo11/lvtest
mke2fs 1.41.9 (22-Aug-2009)
Filesystem label=
OS type: Linux
Block size=4096 (log=2)
Fragment size=4096 (log=2)
130816 inodes, 523264 blocks
26163 blocks (5.00%) reserved for the super user
First data block=0
Maximum filesystem blocks=536870912
16 block groups
32768 blocks per group, 32768 fragments per group
8176 inodes per group
Superblock backups stored on blocks:
    32768, 98304, 163840, 229376, 294912

Writing inode tables: done
Creating journal (8192 blocks): done
Writing superblocks and filesystem accounting information: done

This filesystem will be automatically checked every 21 mounts or
180 days, whichever comes first.  Use tune2fs -c or -i to override.
$
```

After you've created the new filesystem, you can mount the volume in the virtual directory using the standard Linux mount command, just as if it were a physical partition. The only difference is that you use a special path that identifies the logical volume:

```
$ sudo mount /dev/Vol1/lvtest test
$ cd test
$ ls -al
total 24
drwxr-xr-x.  3 root root  4096 2010-09-17 17:36 .
drwx----- 33 rich  rich  4096 2010-09-17 17:37 ..
drwx-----.  2 root root 16384 2010-09-17 17:36 lost+found
$
```

Notice that the paths used in both the `mkfs.ext4` and `mount` commands is a little odd. Instead of a physical partition path, the path uses the volume group name, along with the logical volume name. Once the filesystem is mounted, you can access the new area in the virtual directory.

Modifying the LVM

Because the benefit of using the Linux LVM is to dynamically modify filesystems, you'd expect there to be some tools that allow you to do that. There are some tools available in Linux that allow you to modify the existing Logical Volume Management configuration.

If you don't have access to a fancy graphical interface for managing your Linux LVM environment, all is not lost. You've already seen some of the Linux LVM command line programs in action in this chapter. There are also a host of other command line programs you can use to manage the LVM setup once you've installed it. [Table 7.6](#) lists the common commands that are available in the Linux LVM package.

Table 7.6 The Linux LVM Commands

Command	Function
---------	----------

<code>vgchange</code>	Activate and deactivate a volume group
<code>vgremove</code>	Remove a volume group
<code>vgextend</code>	Add physical volumes to a volume group
<code>vgreduce</code>	Remove physical volumes from a volume group
<code>lvextend</code>	Increase the size of a logical volume
<code>lvreduce</code>	Decrease the size of a logical volume

Using these command line programs, you have full control over your Linux LVM environment.

Warning

Be careful when manually increasing or decreasing the size of a logical volume. The filesystem stored in the logical volume will need to be manually fixed to handle the change in size. Most filesystems include command line programs for reformatting the filesystem, such as the `resize2fs` program for the `ext2` and `ext3` filesystems.

Summary

Working with storage devices in Linux requires that you know a little bit about filesystems. Knowing how to create and work with filesystems from the command line can come in handy as you work on Linux systems. This chapter discussed how to handle filesystems from the Linux command line.

The Linux system is different from Windows in that it

supports lots of different methods for storing files and folders. Each filesystem method has different features that make it ideal for different situations. Also, each filesystem method uses different commands for interacting with the storage device.

Before you can install a filesystem on a storage device, you must first prepare the device. The `fdisk` command is used to partition storage devices to get them ready for the filesystem. When you partition the storage device you must define what type of filesystem will be used on it.

After you partition a storage device, you can use one of several different filesystems for the partition. The most popular Linux filesystems are `ext3` and `ext4`. Both of these filesystems provide journaling filesystem features, making them less prone to errors and problems if the Linux system should crash.

One limiting factor to creating filesystems directly on a storage device partition is that you can't easily change the size of the filesystem if you run out of disk space. However, Linux supports logical volume management, a method of creating virtual partitions across multiple storage devices. This method allows you to easily expand an existing filesystem without having to completely rebuild it. The Linux LVM package provides command line commands to create logical volumes across multiple storage devices to build filesystems on.

Now that you've seen the core Linux command line commands, it's close to time to start creating some shell

script programs. However, before you start coding there's one more element we need to discuss: editors. If you plan on writing shell scripts, you'll need an environment in which to create your masterpieces. The next chapter discusses how to install and manage software packages from the command line in different Linux environments.

Chapter 8

Installing Software

In This Chapter

- Installing software
- Using Debian packages
- Working with Red Hat packages

In the old days of Linux, installing software could be a painful experience. Fortunately, the Linux developers have made life a little easier for us by bundling software into pre-built packages that are much easier to install. However, there's still a little work on our part to get the software packages installed, especially if you want to do that from the command line. This chapter takes a look at the various Package Management Systems available in Linux, and the command line tools used for software installation, management, and removal.

Package Management Primer

Before diving into the world of Linux software package

management, this chapter goes through a few of the basics first. Each of the major Linux distributions utilizes some form of a Package Management System (PMS) to control installing software applications and libraries. A PMS utilizes a database that keeps track of:

- What software packages are installed on the Linux system
- What files have been installed for each package
- Versions of each of the software packages installed

Software packages are stored on servers, called *repositories*, and are accessed across the Internet via PMS utilities running on your local Linux system. You can use the PMS utilities to search for new software packages or even updates to software packages already installed on the system.

A software package will often have *dependencies* or other packages that must be installed first for the software to run properly. The PMS utilities will detect these dependencies and offer to install any additionally needed software packages before installing the desired package.

The downside to PMS is that there isn't a single standard utility. Whereas all of the bash shell commands discussed so far in this book will work no matter which Linux distribution you use, this is not true with software package management.

The PMS utilities and their associated commands are vastly different between the various Linux distributions. The two primary PMS base utilities commonly used in the Linux

world are `dpkg` and `rpm`.

Debian-based distributions such as Ubuntu and Linux Mint use, at the base of their PMS utilities, the `dpkg` command. This command interacts directly with the PMS on the Linux system and is used for installing, managing, and removing software packages.

The Red Hat-based distributions, such as Fedora, openSUSE, and Mandriva, use the `rpm` command at the base of their PMS. Similar to the `dpkg` command, the `rpm` command can list installed packages, install new packages, and remove existing software.

Note that these two commands are the core of their respective PMS, not the entire PMS itself. Many Linux distributions that use the `dpkg` or `rpm` methods have built additional specialty PMS utilities upon these base commands to help make your life much easier. The following sections walk through various PMS utility commands you'll run into in the popular Linux distributions.

The Debian-Based Systems

The `dpkg` command is at the core of the Debian-based family of PMS tools. Other tools included in this PMS are:

- `apt-get`
- `apt-cache`
- `aptitude`

By far the most common command line tool is `aptitude`,

and for good reason. The aptitude tool is essentially a front-end for both the apt tools and `dpkg`. Whereas `dpkg` is a PMS tool, aptitude is a complete Package Management System.

Using the `aptitude` command at the command line will help you avoid common software installation problems, such as missing software dependencies, unstable system environments, and just a whole lot of unnecessary hassle. This section takes a look at how to use the `aptitude` command tool from the Linux command line.

Managing Packages with aptitude

A common task faced by Linux system administrators is to determine what packages are already installed on the system. Fortunately, aptitude has a handy interactive interface that makes this task an easy one.

At the shell prompt, type `aptitude` and press Enter. You will be thrown into aptitude's full-screen mode, as you can see in [Figure 8.1](#).

[Figure 8.1](#) The aptitude main window

```
File Edit View Search Terminal Help
Actions Undo Package Resolver Search Options Views Help
C-T: Menu ?: Help q: Quit u: Update g: Download/Install/Remove Pkgs
aptitude 0.6.3 Will use 4,948kB of disk space DL Size: 1,927k
--- New Packages (6530)
--- Installed Packages (1535)
--- Not Installed Packages (24288)
--- Obsolete and Locally Created Packages (4)
--- Virtual Packages (3568)
--- Tasks (16196)

These packages are currently installed on your computer.

This group contains 1535 packages.
```

Use the arrow keys to maneuver around the menu. Select the menu option Installed Packages to see what packages are installed. You will see several groups of software packages, such as editors, and so on. A number in parentheses follows each group, which indicates the number of packages the group contains.

Use the arrow keys to highlight a group and press Enter to see each subgroup of packages. You will then see the individual package names and their version numbers. Use the Enter key on individual packages to get very detailed information, such as the package's description, home page, size, maintainer, and so on.

When you're done viewing the installed packages, press

q to quit the display. You can then go back to the arrow keys and use Enter to toggle open or closed the packages and their subgroups. When you are all done, just press q multiple times until you receive the pop-up screen “Really quit Aptitude?”

If you already know the packages on your system and want to quickly display detailed information about a particular package, there is no need to go into aptitude's interactive interface. You can use aptitude as a single command at the command line:

```
aptitude show package_name
```

Here's an example of displaying the details of the package grub2-theme-mint:

```
$ aptitude show grub2-theme-mint
Package: grub2-theme-mint
New: yes
State: installed
Automatically installed: no
Version: 1.0.3
Priority: optional
Section: misc
Maintainer: Clement Lefebvre <root@linuxmint.com>
Uncompressed Size: 442k
Description: Grub2 theme for Linux Mint
    Grub2 theme for Linux Mint

$
```

Note

The `aptitude show` command does not indicate that the package is installed on the system. It shows only detailed package information from the software repository.

One detail you cannot get with `aptitude` is a listing of all the files associated with a particular software package. To get this list, you will need to go to the `dpkg` tool itself:

```
dpkg -L package_name
```

Here's an example of using `dpkg` to list all of the files installed as part of the `grub2-theme-mint` package:

```
$
$ dpkg -L grub2-theme-mint
/.
/boot
/boot/boot
/boot/boot/grub
/boot/boot/grub/linuxmint.png
/boot/grub
/boot/grub/linuxmint.png
/usr
/usr/share
/usr/share/doc
/usr/share/doc/grub2-theme-mint
/usr/share/doc/grub2-theme-mint/changelog.gz
/usr/share/doc/grub2-theme-mint/copyright
/etc
/etc/grub.d
/etc/grub.d/06_mint_theme
$
```

You can also do the reverse—find what package a particular file belongs to:

```
dpkg --search absolute_file_name
```

Note that you need to use an absolute file reference for this to work:

```
$
$ dpkg --search /boot/grub/linuxmint.png
grub2-theme-mint: /boot/grub/linuxmint.png
$
```

The output shows the [linuxmint.png](#) file was installed as part of the `grub2-theme-mint` package.

Installing Software Packages with aptitude

Now that you know more about listing software package information on your system, this section walks through a software package install. First, you'll want to determine the package name to install. How do you find a particular software package? Use the `aptitude` command with the search option:

```
aptitude search package_name
```

The beauty of the search option is that you do not need to insert wildcards around `package_name`. Wildcards are implied. Here's an example of using `aptitude` to look for the wine software package:

```
$
$ aptitude search wine
p  gnome-wine-icon-theme          - red variation of the GNOME- ...
v  libkwineffectsl-api            -
p  libkwineffectsla               - library used by effects...
p  q4wine                         - Qt4 GUI for wine (W.I.N.E)
p  shiki-wine-theme               - red variation of the Shiki- ...
p  wine                           - Microsoft Windows Compatibility ...
p  wine-dev                       - Microsoft Windows Compatibility ...
p  wine-gecko                     - Microsoft Windows Compatibility ...
p  winel.0                        - Microsoft Windows Compatibility ...
p  winel.0-dev                    - Microsoft Windows Compatibility ...
p  winel.0-gecko                  - Microsoft Windows Compatibility ...
p  winel.2                        - Microsoft Windows Compatibility ...
p  winel.2-dbg                    - Microsoft Windows Compatibility ...
p  winel.2-dev                    - Microsoft Windows Compatibility ...
p  winel.2-gecko                  - Microsoft Windows Compatibility ...
p  winefish                       - LaTeX Editor based on Bluefish
$
```

Notice that before each package name is either a `p` or an `i`. If you see an `i` that means that the package is currently installed on your system. If you see a `p` that means it is

available but not installed. As you can see from the preceding listing, this system does not have wine currently installed on it, but the package is available from the software repository.

Installing a software package on a system from a repository using `aptitude` is as easy as this:

```
aptitude install package_name
```

Once you find the software package name from the search option, just plug it into the `aptitude` command using the install option:

```
$
$ sudo aptitude install wine
The following NEW packages will be installed:
  cabextract{a} esound-clients{a} esound-common{a} gnome-exe-thumbnailer
{a}
  icoutils{a} imagemagick{a} libaudio2{a} libaudiofile0{a} libcdt4{a}
  libesd0{a} libgraph4{a} libgvc5{a} libilmbase6{a} libmagickcore3-extra
{a}
  libmpg123-0{a} libnetpbm10{a} libopenall{a} libopenexr6{a}
  libpathplan4{a} libxdot4{a} netpbm{a} ttf-mscorefonts-installer{a}
  ttf-symbol-replacement{a} winbind{a} wine winel.2{a} winel.2-gecko{a}
0 packages upgraded, 27 newly installed, 0 to remove and 0 not upgraded.
Need to get 0B/27.6MB of archives. After unpacking 121MB will be used.
Do you want to continue? [Y/n/?] Y
Preconfiguring packages ...
...
All done, no errors.
All fonts downloaded and installed.
Updating fontconfig cache for /usr/share/fonts/truetype/msttcorefonts
Setting up winbind (2:3.5.4~dfsg-lubuntu7) ...
  * Starting the Winbind daemon winbind
    [ OK ]
Setting up wine (1.2-0ubuntu5) ...
Setting up gnome-exe-thumbnailer (0.6-0ubuntu1) ...
Processing triggers for libc-bin ...
ldconfig deferred processing now taking place

$
```

Note

Before the `aptitude` command in the preceding listing, the `sudo` command is used. The `sudo` command allows you to run a command as the root user. You can use the `sudo` command to run administrative tasks, such as installing software.

To check if the installation processed properly, just use the `search` option again. This time you should see an `i u` listed in front of the wine software package, indicating it is installed.

You may also notice that there are additional packages with the `i u` in front of them. This is because `aptitude` automatically resolved any necessary package dependencies for us and installs the needed additional library and software packages. This is a wonderful feature included in many Package Management Systems.

Updating Software with aptitude

While `aptitude` helps protect you from problems installing software, trying to coordinate a multiple-package update with dependencies can get tricky. To safely update all of the software packages on a system with any new versions in the repository, use the `safe-upgrade` option:

```
aptitude safe-upgrade
```

Notice that this command doesn't take a software package name as an argument. That's because the `safe-upgrade` option will upgrade all the installed packages to the

most recent version available in the repository, which is safer for system stabilization.

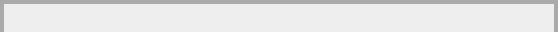
Here's a sample output from running the `aptitude safe-update` command:

```
$
$ sudo aptitude safe-upgrade
The following packages will be upgraded:
  evolution evolution-common evolution-plugins gsfontr libevolutio
  xserver-xorg-video-geode
6 packages upgraded, 0 newly installed, 0 to remove and 0 not upgraded.
Need to get 9,312kB of archives. After unpacking 0B will be used.
Do you want to continue? [Y/n/?] Y
Get:1 http://us.archive.ubuntu.com/ubuntu/ maverick/main
  libevolutio i386 2.30.3-lubuntu4 [2,096kB]
...
Preparing to replace xserver-xorg-video-geode 2.11.9-2
(using .../xserver-xorg-video-geode_2.11.9-3_i386.deb) ...
Unpacking replacement xserver-xorg-video-geode ...
Processing triggers for man-db ...
Processing triggers for desktop-file-utils ...
Processing triggers for python-gmenu ...
...
Current status: 0 updates [-6].
$
```

There are also less conservative options for software upgrades that you can use:

- `aptitude full-upgrade`
- `aptitude dist-upgrade`

These options perform the same task, upgrading all of the software packages to the latest versions. Where they differ from `safe-upgrade` is that they do not check dependencies between packages. The whole package dependency issue can get real ugly. If you're not exactly sure of the dependencies for the various packages, then stick with the `safe-upgrade` option.



Note

Obviously, running aptitude's `safe-upgrade` option is something you should do on a regular basis to keep your system up-to-date. However, it is especially important to run it after a fresh distribution installation. Usually there are lots of security patches and updates that are released since the last full release of a distribution.

Uninstalling Software with aptitude

Getting rid of software packages with aptitude is as easy as installing and upgrading them. The only real choice you have to make is whether or not to keep the software's data and configuration files around afterwards.

To remove a software package, but not the data and configuration files, use the `remove` option of aptitude. To remove a software package and the related data and configuration files, use the `purge` option:

```
$ sudo aptitude purge wine
```

```
[sudo] password for user:
```

```
The following packages will be REMOVED:
```

```
 cabextract{u} esound-clients{u} esound-common{u} gnome-exe-thumbnailer
{u}
 icoutils{u} imagemagick{u} libaudio2{u} libaudiofile0{u} libcdt4{u}
 libesd0{u} libgraph4{u} libgvc5{u} libilmbase6{u} libmagickcore3-extra
{u}
 libmpgl23-0{u} libnetpbm10{u} libopenall{u} libopenexr6{u}
 libpathplan4{u} libxdot4{u} netpbm{u} ttf-mscorefonts-installer{u}
 ttf-symbol-replacement{u} winbind{u} wine{p} winel.2{u} winel.2-gecko
{u}
```

```
0 packages upgraded, 0 newly installed, 27 to remove and 6 not upgraded.
Need to get 0B of archives. After unpacking 121MB will be freed.
```

```
Do you want to continue? [Y/n/?] Y
```

```
(Reading database ... 120968 files and directories currently installed.)
Removing ttf-mscorefonts-installer ...
```

```
...
```

```
Processing triggers for fontconfig ...  
Processing triggers for ureadahead ...  
Processing triggers for python-support ...
```

\$

To see if the package has been removed, you can use the `aptitude search` option again. If you see a `c` in front of the package name, it means the software has been removed, but the configuration files have not been purged from the system. A `p` in front indicates the configuration files have also been removed.

The aptitude Repositories

The default software repository locations for `aptitude` are set up for you when you install your Linux distribution. The repository locations are stored in the file `/etc/apt/sources.list`.

In many cases, you will never need to add/remove a software repository so you won't need to touch this file. However, `aptitude` will only pull software from these repositories. Also, when searching for software to install or update, `aptitude` will only check these repositories. If you need to include some additional software repositories for your PMS, this is the place to do it.

Tip

The Linux distribution developers work hard to make sure package versions added to the repositories don't conflict with one another. Usually it's safest to upgrade or install a

software package from the repository. Even if a newer version is available elsewhere, you may want to hold off installing it until that version is available in your Linux distribution's repository.

The following is an example of a `sources.list` file from an Ubuntu system:

```
$
$ cat /etc/apt/sources.list
#deb cdrom:[Ubuntu 10.10 _Maverick Meerkat_ - Alpha i386
(20100921.1)]/ maverick main restricted
# See http://help.ubuntu.com/community/UpgradeNotes for how to upgrade
to
# newer versions of the distribution.

deb http://us.archive.ubuntu.com/ubuntu/ maverick main restricted
deb-src http://us.archive.ubuntu.com/ubuntu/ maverick main restricted

## Major bug fix updates produced after the final release of the
## distribution.
deb http://us.archive.ubuntu.com/ubuntu/ maverick-updates main
restricted
deb-src http://us.archive.ubuntu.com/ubuntu/ maverick-updates main
restricted
...
## This software is not part of Ubuntu, but is offered by third-party
## developers who want to ship their latest software.
deb http://extras.ubuntu.com/ubuntu maverick main
deb-src http://extras.ubuntu.com/ubuntu maverick main

deb http://security.ubuntu.com/ubuntu maverick-security main restricted
deb-src http://security.ubuntu.com/ubuntu maverick-security main
restricted
deb http://security.ubuntu.com/ubuntu maverick-security universe
deb-src http://security.ubuntu.com/ubuntu maverick-security universe
deb http://security.ubuntu.com/ubuntu maverick-security multiverse
deb-src http://security.ubuntu.com/ubuntu maverick-security multiverse
$
```

First, notice that the file is full of helpful comments and warnings. The repository sources specified use the following structure:

```
deb (or deb-src) address distribution_name package_type_list
```

The `deb` or `deb-src` value indicates the software package type. The `deb` value indicates it is a source of compiled programs, whereas the `deb-src` value indicates it is a source of source code.

The `address` entry is the software repository's Web address. The `distribution_name` entry is the name of this particular software repository's distribution's version. In the example, the distribution name is `maverick`. This does not necessarily mean that the distribution you are running is Ubuntu's Maverick Meercat; it just means the Linux distribution is using the Ubuntu Maverick Meercat software repositories! For example, in Linux Mint's `sources.list` file, you will see a mix of Linux Mint and Ubuntu software repositories.

Finally, the `package_type_list` entry may be more than one word and indicates what type of packages the repository has in it. For example, you may see values such as `main`, `restricted`, `universe`, or `partner`.

When you need to add a software repository to your `sources` file, you can try to wing it yourself, but that more than likely will cause problems. Often, software repository sites or various package developer sites will have an exact line of text that you can copy from their Website and paste into your `sources.list` file. It's best to choose the safer route and just copy/paste.

The front-end interface, `aptitude`, provides intelligent command line options for working with the Debian-based `dpkg` utility. Now it's time to take a look at the Red Hat–

based distributions' `rpm` utility and its various front-end interfaces.

The Red Hat–Based Systems

Like the Debian-based distributions, the Red Hat–based systems have several different front-end tools that are available. The common ones are:

- `yum`: Used in Red Hat and Fedora
- `urpm`: Used in Mandriva
- `zypper`: Used in openSUSE

These front-ends are all based on the `rpm` command line tool. The following section discusses how to manage software packages using these various `rpm`-based tools. The focus will be on `yum`, but information will also be included for `zypper` and `urpm`.

Listing Installed Packages

To find out what is currently installed on your system, at the shell prompt type the following command:

```
yum list installed
```

The information will probably whiz by you on the display screen, so it's best to redirect the installed software listing into a file. You can then use the `more` or `less` command (or a GUI editor) to look at the list in a controlled manner.

```
yum list installed > installed_software
```


To list out the installed packages on your openSUSE or Mandriva distribution, see the commands in [Table 8.1](#). Unfortunately, the `urpm` tool used in Mandriva cannot produce a currently installed software listing. Thus, you will need to revert to the underlying `rpm` tool.

Table 8.1 How to List Installed Software with zypper and urpm

Distribution	Front-End Tool	Command
Mandriva	urpm	<code>rpm -qa > installed_software</code>
openSUSE	zypper	<code>zypper search -I > installed_software</code>

To find out detailed information for a particular software package, `yum` really shines. Not only will it give you a very verbose description of the package, but with another simple command, you can see if the package is installed:

```
# yum list xterm
Loaded plugins: langpacks, presto, refresh-packagekit
Adding en_US to language list
Available Packages
xterm.i686 261-2.fc14 fedora
#
# yum list installed xterm
Loaded plugins: refresh-packagekit
Error: No matching Packages to list
#
```

The commands to list detailed software package information using `urpm` and `zypper` are in [Table 8.2](#). You can acquire an even more detailed set of package information from the repository, using the `info` option on the `zypper` command.

Table 8.2 How to See Various Package Details with

zypper and urpm

Detail Type	Front-End Tool	Command
Package Information	urpm	<code>urpmq -i <i>package_name</i></code>
Installed?	urpm	<code>rpm -q <i>package_name</i></code>
Package Information	zypper	<code>zypper search -s <i>package_name</i></code>
Installed?	zypper	Same command, but look for an <code>i</code> in the Status column

Finally, if you need to find out what software package provides a particular file on your filesystem, the versatile `yum` can do that, too! Just enter the command:

```
yum provides file_name
```

Here's an example of trying to find what software provided the configuration file `/etc/yum.conf`:

```
#
# yum provides /etc/yum.conf
Loaded plugins: langpacks, presto, refresh-packagekit
Adding en_US to language list
yum-3.2.28-5.fc14.noarch : RPM installer/updater
Repo : fedora
Matched from:
Filename : /etc/yum.conf

yum-3.2.28-5.fc14.noarch : RPM installer/updater
Repo : installed
Matched from:
Other : Provides-match: /etc/yum.conf

#
```

`yum` checked two separate repositories: `fedora` and `installed`. From both, the answer is: the `yum` software package provides this file!

Installing Software with yum

Installation of a software package using `yum` is incredibly easy. The following is the basic command for installing a software package, all its needed libraries, and package dependencies from a repository:

```
yum install package_name
```

Here's an example of installing the `xterm` package:

```
$ su -
Password:
# yum install xterm
Loaded plugins: langpacks, presto, refresh-packagekit
Adding en_US to language list
fedora/metalink | 20 kB 00:00
fedora | 4.3 kB 00:00
fedora/primary_db | 11 MB 01:57
updates/metalink | 16 kB 00:00
updates | 4.7 kB 00:00
updates/primary_db | 3.1 MB 00:30
Setting up Install Process
Resolving Dependencies
--> Running transaction check
--> Package xterm.i686 0:261-2.fc14 set to be installed
...
Installed:
xterm.i686 0:261-2.fc14

Complete!
#
#
```

Note

Before the `yum` command in the preceding listing, the `su -` command is used. This command allows you to switch to the root user. On this Linux system, the `#` denotes you are logged in as root. You should only switch to root user temporarily in order to run

administrative tasks, such as installing and updating software. The `sudo` command is another option as well.

You can also manually download an `rpm` installation file and install it using `yum`. This is called a *local installation*. The basic command is:

```
yum localinstall package_name.rpm
```

You can begin to see that one of `yum`'s strengths is that it uses very logical and user-friendly commands.

[Table 8.3](#) shows how to perform a package install with `urpm` and `zypper`. You should note that if you are not logged in as root, you will get a “command not found” error message using `urpm`.

[Table 8.3](#) How to Install Software with `zypper` and `urpm`

Front-End Tool	Command
<code>urpm</code>	<code>urpmi package_name</code>
<code>zypper</code>	<code>zypper install package_name</code>

Updating Software with `yum`

In most Linux distributions, when you're working away in the GUI, you get those nice little notification icons telling you an update is needed. Here at the command line, it takes a little more work.

To see the list of all the available updates for your installed packages, type the following command:

```
yum list updates
```

It's always nice to get no response to this command because it means you have nothing to update! However, if you do discover a particular software package needs updating, then type in the following command:

```
yum update package_name
```

If you'd like to update all the packages listed in the update list, just enter the following command:

```
yum update
```

Commands for updating software packages on Mandriva and openSUSE are listed in [Table 8.4](#). When urpm is used, the repository database will be automatically refreshed as well as software packages updated.

Table 8.4 How to Update Software with zypper & urpm

Front-End Tool	Command
urpm	urpmi --auto-update --update
zypper	zypper update

Uninstalling Software with yum

The `yum` tool also provides an easy way to uninstall software you no longer want on your system. As with `aptitude`, you need to choose whether or not to keep the software package's data and configuration files.

To just remove the software package and keep any configuration and data files, use the following command:

```
yum remove package_name
```

To uninstall the software and all its files, use the `erase`

option:

```
yum erase package_name
```

You will see it is equally easy to remove software using `urpm` and `zypper` in [Table 8.5](#). Both of these tools perform a function similar to `yum`'s `erase` option.

Table 8.5 How to Uninstall Software with `zypper` and `urpm`

Front-End Tool	Command
<code>urpm</code>	<code>urpme package_name</code>
<code>zypper</code>	<code>zypper remove package_name</code>

While life is considerably easier with PMS packages, it's not always problem free. Occasionally things do go wrong. Fortunately, there's help.

Dealing with Broken Dependencies

Sometimes as multiple software packages get loaded, a software dependency for one package can get overwritten by the installation of another package. This is called a *broken dependency*.

If this should happen on your system, first try the following command:

```
yum clean all
```

Then try to use the `update` option in the `yum` command. Sometimes, just cleaning up any misplaced files can help.

If that doesn't solve the problem, try the following command:

```
yum deplist package_name
```

This command displays all the package's library dependencies and what software package provides them. Once you know the libraries required for a package, you can then install them. Here's an example of determining the dependencies for the xterm package:

```
# yum deplist xterm
Loaded plugins: langpacks, presto, refresh-packagekit
Adding en_US to language list
Finding dependencies:
package: xterm.i686 261-2.fc14
dependency: libutempter.so.0
provider: libutempter.i686 1.1.5-4.fc12
dependency: rtld(GNU_HASH)
provider: glibc.i686 2.12.90-17
provider: glibc.i686 2.12.90-21
dependency: libc.so.6(GLIBC_2.4)
provider: glibc.i686 2.12.90-17
provider: glibc.i686 2.12.90-21
...
dependency: /bin/sh
provider: bash.i686 4.1.7-3.fc14
dependency: libICE.so.6
provider: libICE.i686 1.0.6-2.fc13
dependency: libXmu.so.6
provider: libXmu.i686 1.0.5-2.fc13
dependency: libc.so.6(GLIBC_2.3)
provider: glibc.i686 2.12.90-17
provider: glibc.i686 2.12.90-21
dependency: libXaw.so.7
provider: libXaw.i686 1.0.6-4.fc12
dependency: libX11.so.6
provider: libX11.i686 1.3.4-3.fc14
dependency: libc.so.6(GLIBC_2.2)
provider: glibc.i686 2.12.90-17
provider: glibc.i686 2.12.90-21
#
```

If that doesn't solve your problem, you have one last tool:

```
yum update --skip-broken
```

The `--skip-broken` option allows you to just ignore the package with the broken dependency and update the other software packages. This may not help the broken package,

but at least you can update the remaining packages on the system!

In [Table 8.6](#), the commands to try for broken dependencies with `urpm` and `zypper` are listed. With `zypper`, there is only the one command to verify and fix a broken dependency. With `urpm`, if the `clean` option does not work, you can skip updates on the offensive package. To do this, you must add the name of the offending package to the file `/etc/urpmi/skip.list`.

Table 8.6 Broken Dependencies with `zypper` and `urpm`

Front End Tool	Command
<code>urpm</code>	<code>urpmi --clean</code>
<code>zypper</code>	<code>zypper verify</code>

Yum Repositories

Just like the `aptitude` systems, `yum` has its software repositories set up at installation. For most purposes, these pre-installed repositories will work just fine for your needs. But if and when the time comes that you need to install software from a different repository, here are some things you will need to know.

Tip

A wise system administrator sticks with approved repositories. An approved repository is one that is sanctioned by the distribution's official site. If you start adding

unapproved repositories, you lose the guarantee of stability.
And you will be heading into broken dependencies territory.

To see what repositories you are currently pulling software from, type in the following command:

```
yum repolist
```

If you don't find a repository you need software from, then you will need to do a little configuration file editing. The `yum` repository definition files are located in `/etc/yum.repos.d`. You will need to add the proper URL and gain access to any necessary encryption keys.

Good repository sites such as rpmfusion.org will lay out all the steps necessary to use them. Sometimes these repository sites will offer an `rpm` file that you can download and install using the `yum localinstall` command. The installation of the `rpm` file will do all the repository setup work for you. Now that's convenient!

`urpm` calls its repositories *media*. The commands for looking at `urpm` media and `zypper`'s repositories are in [Table 8.7](#). You will notice with both of these front-end tools that you do not edit a configuration file. Instead, to add media or a repository, you just type in the command.

[Table 8.7](#) `zypper` and `urpm` Repositories

Action	Front-End Tool	Command
Display repository	<code>urpm</code>	<code>urpmq --list-media</code>
Add repository	<code>urpm</code>	<code>urpmi.addmedia path_name</code>
Display repository	<code>zypper</code>	<code>zypper repos</code>

Both Debian- and Red Hat-based systems use Package Management Systems to ease the process of managing software. Now we are going to step out of the world of Package Management Systems and look at something a little more difficult, installing directly from source code.

Installing from Source Code

Chapter 4 discussed *tarball* packages—how to create them using the `tar` command line command and how to unpack them. Before the fancy `rpm` and `dpkg` tools, administrators had to know how to unpack and install software from tarballs.

If you work in the open source software environment much, there's a good chance you will still find software packed up as a tarball. This section walks you through the process of unpacking and installing a tarball software package.

For this example, the software package `sysstat` will be used. The `sysstat` utility is a very nice software package that provides a variety of system monitoring tools.

First, you will need to download the `sysstat` tarball to your Linux system. While you can often find the `sysstat` package available on different Linux sites, it's usually best to go straight to the source of the program. In this case, it's the

Website <http://sebastien.godard.pagesperso-orange.fr/>.

If you click the Download link, you'll go to the page that contains the files for downloading. The current version at the time of this writing is 9.1.5, and the distribution file name is `sysstat-9.1.5.tar.gz`.

Click the link to download the file to your Linux system. Once you have downloaded the file, you can unpack it.

To unpack a software tarball, use the standard `tar` command:

```
#
# tar -zxvf sysstat-9.1.5.tar.gz
sysstat-9.1.5/
sysstat-9.1.5/sar.c
sysstat-9.1.5/iostat.c
sysstat-9.1.5/sadc.c
sysstat-9.1.5/sa.h
sysstat-9.1.5/iconfig
sysstat-9.1.5/CHANGES
sysstat-9.1.5/COPYING
sysstat-9.1.5/CREDITS
sysstat-9.1.5/sa2.in
sysstat-9.1.5/README
sysstat-9.1.5/crontab.sample
sysstat-9.1.5/nls/
...
sysstat-9.1.5/nfsiostat.c
sysstat-9.1.5/sysstat-9.1.5.lsm
sysstat-9.1.5/cifsioestat.c
sysstat-9.1.5/nfsiostat.h
sysstat-9.1.5/cifsioestat.h
sysstat-9.1.5/sysstat-9.1.5.spec
#
```

Now that the tarball is unpacked and the files have neatly put themselves into a directory called `sysstat-9.1.5`, you can dive down into that directory and continue.

First, use the `cd` command to get into the new directory and then list the contents of the directory:

```
$ cd sysstat-9.1.5
```

```

$ ls
activity.c      INSTALL      prf_stats.h  sysconfig.in
build          ioconf.c    pr_stats.c   sysstat-9.1.5.lsm
CHANGES       ioconf.h    pr_stats.h   sysstat-9.1.5.spec
cifsiostat.c   iostat.c    rd_stats.c   sysstat.cron.daily.in
cifsiostat.h   iostat.h    rd_stats.h   sysstat.cron.in
common.c       Makefile.in README       sysstat.cron.hourly.in
common.h       man         sal.in       sysstat.in
configure      mpstat.c    sa2.in       sysstat.ioconf
configure.in    mpstat.h    sa_common.c  sysstat.sysconfig.in
contrib        nfsiostat.c sadc.c       TODO
COPYING        nfsiostat.h sadf.c       version.in
CREDITS        nls         sadf.h       xml
crontab.sample pidstat.c   sa.h
FAQ            pidstat.h   sar.c
iconfig        prf_stats.c sa_wrap.c
$

```

In the listing out of the directory, you should typically see a `README` or `AAAREADME` file. It is very important to read this file. In the file will be the actual instructions you will need to finish the software's installation.

Following the advice contained in the `README` file, the next step is to `configure sysstat` for your system. This checks your Linux system to ensure it has the proper library dependencies, in addition to the proper compiler to compile the source code:

```

# ./configure
Check programs:
.
checking for gcc... gcc
checking for C compiler default output file name... a.out
checking whether the C compiler works... yes
...
checking for ANSI C header files... (cached) yes
checking for dirent.h that defines DIR... yes
checking for library containing opendir... none required
checking ctype.h usability... yes
checking ctype.h presence... yes
checking for ctype.h... yes
checking errno.h usability... yes
...
Check library functions:

```

```
.
checking for strchr... yes
checking for strcspn... yes
checking for strspn... yes
checking for strstr... yes
checking for sensors support... yes
...
Check configuration:
...
config.status: creating Makefile

Sysstat version:          9.1.5
Installation prefix:      /usr/local
rc directory:             /etc/rc.d
Init directory:          /etc/rc.d/init.d
Configuration directory: /etc/sysconfig
Man pages directory:     /usr/local/man
Compiler:                 gcc
Compiler flags:           -g -O2
```

```
#
```

If anything does go wrong, the `configure` step will display an error message explaining what's missing.

The next stage is to build the various binary files using the `make` command. The `make` command compiles the source code and then the linker to create the final executable files for the package. As with the `configure` command, the `make` command produces lots of output as it goes through the steps of compiling and linking all of the source code files:

```
# make
...
gcc -o nfsiostat -g -O2 -Wall -Wstrict-prototypes -pipe
-O2 nfsiostat.o librdstats.a libsyscom.a -s
gcc -o cfsiostat.o -c -g -O2 -Wall -Wstrict-prototypes -pipe
-O2 -DSA_DIR="/var/log/sa/"
-DSADC_PATH="/usr/local/lib/sa/sadc/" cfsiostat.o
gcc -o cfsiostat -g -O2 -Wall -Wstrict-prototypes -pipe
-O2 cfsiostat.o librdstats.a libsyscom.a -s
#
```

When `make` is finished, you'll have the actual `sysstat` software program available in the directory! However, it's

somewhat inconvenient to have to run it from that directory. Instead, you'll want to install it in a common location on your Linux system. To do that, you'll need to log in as the root user account (or use the `sudo` command if your Linux distribution prefers), and then use the `install` option of the `make` command:

```
# make install
mkdir -p /usr/local/man/man1
mkdir -p /usr/local/man/man8
rm -f /usr/local/man/man8/sal.8*
install -m 644 -g man man/sal.8 /usr/local/man/man8
rm -f /usr/local/man/man8/sa2.8*
install -m 644 -g man man/sa2.8 /usr/local/man/man8
rm -f /usr/local/man/man8/sadc.8*
install -m 644 -g man man/sadc.8 /usr/local/man/man8
rm -f /usr/local/man/man1/sar.1*
...
install -m 644 sysstat.sysconfig /etc/sysconfig/sysstat
install -m 644 CHANGES /usr/local/share/doc/sysstat-9.1.5
install -m 644 COPYING /usr/local/share/doc/sysstat-9.1.5
install -m 644 CREDITS /usr/local/share/doc/sysstat-9.1.5
install -m 644 README /usr/local/share/doc/sysstat-9.1.5
install -m 644 FAQ /usr/local/share/doc/sysstat-9.1.5
install -m 644 *.lsm /usr/local/share/doc/sysstat-9.1.5
#
```

Now the `sysstat` package is installed on the system! While not quite as easy as installing a software package via a PMS, installing software using tarballs is not that difficult.

Summary

This chapter discussed how to work with a software Package Management Systems (PMS) to install, update, or remove software from the command line. While most of the Linux distributions use fancy GUI tools for software

package management, you can also perform package management from the command line.

The Debian-based Linux distributions use the `dpkg` utility to interface with the PMS from the command line. A front-end to the `dpkg` utility is `aptitude`. It provides simple command line options for working with software packages in the `dpkg` format.

The Red Hat-based Linux distributions are based on the `rpm` utility but use different front-end tools at the command line. Red Hat and Fedora use `yum` for installing and managing software packages. The openSUSE distribution uses `zypper` for managing software, while the Mandriva distribution uses `urpm`.

The chapter closed with a discussion on how to install software packages that are only distributed in source code tarballs. The `tar` command allows you to unpack the source code files from the tarball, and then `configure` and `make` allow you to build the final executable program from the source code.

The next chapter takes a look at the different editors available in Linux distributions. As you get ready to start working on shell scripts, it will come in handy to know what editors are available to use!

Chapter 9

Working with Editors

In This Chapter

- Working with the vim editor
- Understanding emacs
- Getting comfortable with KDE
- Using the GNOME editor

Before you can start your shell scripting career, you'll need to know how to use at least one text editor in Linux. The more you know about how to use these fancy features such as searching, cutting, and pasting, the quicker you'll be able to develop your shell scripts. This chapter discusses the main text editors you'll see in the Linux world.

The vim Editor

If you're working in command line mode, you may want to become familiar with at least one text editor that operates in the Linux console. The vi editor is the original editor used on Unix systems. It uses the console graphics mode to

emulate a text-editing window, allowing you to see the lines of your file, move around within the file, and insert, edit, and replace text.

While it is quite possibly the most complicated editor in the world (at least in the opinion of those who hate it), it provides many features that have made it a staple for Unix administrators for decades.

When the GNU Project ported the `vi` editor to the open source world, they chose to make some improvements to it. Because it no longer resembled the original `vi` editor found in the Unix world, the developers also renamed it, to `vi` improved, or *vim*.

To make life easier, almost all Linux distributions create an alias named (see Chapter 5) `vi` that points to the `vim` program:

```
$ alias vi
alias vi='vim'
$
```

This section walks you through the basics of using the `vim` editor to edit your text shell script files.

The Basics of vim

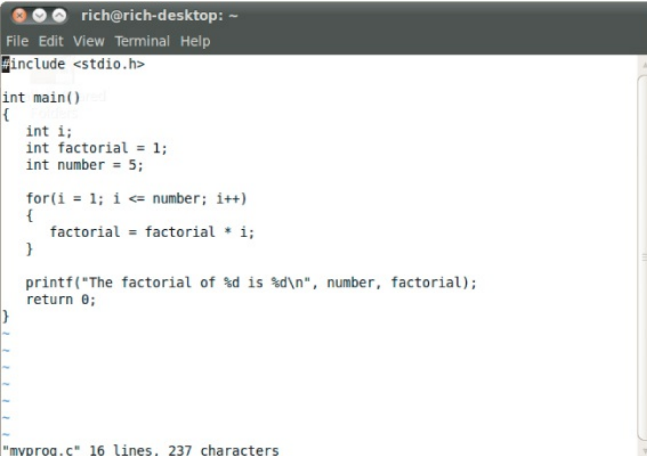
The `vim` editor works with data in a memory buffer. To start the `vim` editor, just type the `vim` command (or `vi` if there's an alias) and the name of the file you want to edit:

```
$ vim myprog.c
```

If you start `vim` without a filename, or if the file doesn't exist, `vim` opens a new buffer area for editing. If you specify

an existing file on the command line, vim will read the entire contents of the file into a buffer area, where it is ready for editing, as shown in [Figure 9.1](#).

Figure 9.1 The vim main window

A screenshot of a terminal window titled "rich@rich-desktop: ~". The terminal shows the vim editor interface with a menu bar (File, Edit, View, Terminal, Help) and a text area containing C code. The code defines a function main() that calculates the factorial of 5. At the bottom, a status line indicates the file "myprog.c" is 16 lines and 237 characters long. The text area has a scrollbar on the right, and the bottom of the window shows several tilde (~) characters representing lines not currently visible.

```
rich@rich-desktop: ~
File Edit View Terminal Help

#include <stdio.h>

int main()
{
    int i;
    int factorial = 1;
    int number = 5;

    for(i = 1; i <= number; i++)
    {
        factorial = factorial * i;
    }

    printf("The factorial of %d is %d\n", number, factorial);
    return 0;
}

~
~
~
~
~

"myprog.c" 16 lines, 237 characters
```

The vim editor detects the terminal type for the session (see Chapter 2) and uses a full-screen mode to use the entire console window for the editor area.

The initial vim edit window shows the contents of the file (if there are any) along with a message line at the bottom of the window. If the file contents don't take up the entire screen, vim places a tilde on lines that are not part of the file (as shown in [Figure 9.1](#)).

The message line at the bottom indicates information about the edited file, depending on the status of the file, and the default settings in your vim installation. If the file is new, the message `[New File]` appears.

The vim editor has two modes of operation:

- Normal mode
- Insert mode

When you first open a file (or start a new file) for editing, the vim editor enters *normal* mode. In normal mode, the vim editor interprets keystrokes as commands (more on those later).

In *insert* mode, vim inserts every key you type at the current cursor location in the buffer. To enter insert mode, press the `i` key. To get out of insert mode and go back into normal mode, press the Escape key on the keyboard.

In normal mode, you can move the cursor around the text area by using the arrow keys (as long as your terminal type is detected properly by vim). If you happen to be on a flaky terminal connection that doesn't have the arrow keys defined, all hope is not lost. The vim commands include commands for moving the cursor:

- `h` to move left one character
- `j` to move down one line (the next line in the text)
- `k` to move up one line (the previous line in the text)
- `l` to move right one character

Moving around within large text files line by line can get tedious. Fortunately, vim provides a few commands to help speed things along:

- `PageDown` (or `Ctrl+F`) to move forward one screen of data
- `PageUp` (or `Ctrl+B`) to move backward one screen of data
- `G` to move to the last line in the buffer
- `num G` to move to the line number `num` in the buffer
- `gg` to move to the first line in the buffer

The vim editor has a special feature within normal mode called *command line mode*. The command line mode provides an interactive command line where you can enter additional commands to control the actions in vim. To get to command line mode, press the colon key in normal mode. The cursor moves to the message line, and a colon appears, waiting for you to enter a command.

Within the command line mode are several commands for saving the buffer to the file and exiting vim:

- `q` to quit if no changes have been made to the buffer data
- `q!` to quit and discard any changes made to the buffer data
- `w filename` to save the file under a different filename
- `wq` to save the buffer data to the file and quit

After seeing just a few basic vim commands you might understand why some people absolutely hate the vim editor. To be able to use vim to its fullest, you must know plenty of obscure commands. However, once you get a few of the basic vim commands down, you can quickly edit files

directly from the command line, no matter what type of environment you're in. Plus, once you get comfortable typing commands, it almost seems second-nature to type both data and editing commands, and it becomes odd having to jump back to using a mouse!

Editing Data

While in insert mode, you can insert data into the buffer; however, sometimes you need to add or remove data after you've already entered it into the buffer. While in normal mode, the vim editor provides several commands for editing the data in the buffer. [Table 9.1](#) lists some common editing commands for vim.

Table 9.1 vim Editing Commands

Command	Description
x	Delete the character at the current cursor position.
dd	Delete the line at the current cursor position.
dw	Delete the word at the current cursor position.
d\$	Delete to the end of the line from the current cursor position.
J	Delete the line break at the end of the line at the current cursor position (joins lines).
u	Undo the previous edit command.
a	Append data after the current cursor position.
A	Append data to the end of the line at the current cursor position.
r char	Replace a single character at the current cursor position with char.

Some of the editing commands also allow you to use a numeric modifier to indicate how many times to perform the command. For example, the command `2x` deletes two characters, starting from the current cursor position, and the command `5dd` deletes five lines, starting at the line from the current cursor position.

Caution

Be careful when trying to use the PC keyboard Backspace or Delete keys while in the vim editor. The vim editor usually recognizes the Delete key as the functionality of the `x` command, deleting the character at the current cursor location. Usually, the vim editor doesn't recognize the Backspace key.

Copy and Paste

A standard feature of modern editors is the ability to cut or copy data, then paste it elsewhere in the document. The vim editor provides a way to do this.

Cutting and pasting is relatively easy. You've already seen the commands in [Table 9.1](#) that can remove data from the buffer. However, when vim removes data, it actually keeps it stored in a separate register. You can retrieve that data by using the `p` command.

For example, you can use the `dd` command to delete a

line of text, then move the cursor to the location in the buffer where you want to place it, and then use the `p` command. The `p` command inserts the text after the line at the current cursor position. You can do this with any command that removes text.

Copying text is a little bit trickier. The copy command in vim is `y` (for yank). You can use the same second character with `y` as with the `d` command (`yw` to yank a word, `y$` to yank to the end of a line). After you yank the text, move the cursor to the location where you want to place the text, and use the `p` command. The yanked text now appears at that location.

Yanking is tricky in that you can't see what happened because you're not affecting the text that you yank. You never know for sure what you yanked until you paste it somewhere. But there's another feature in vim that helps you out with yanking.

The *visual mode* highlights text as you move the cursor. You use visual mode to select text to yank for pasting. To enter visual mode, move the cursor to the location where you want to start yanking, and press `v`. You'll notice that the text at the cursor position is now highlighted. Next, move the cursor to cover the text you want to yank (you can even move down lines to yank more than one line of text). As you move the cursor, vim highlights the text in the yank area. After you've covered the text you want to copy, press the `y` key to activate the yank command. Now that you've got the text in the register, just move the cursor to where you want to paste, and use the `p` command.

Search and Substitute

You can easily search for data in the buffer using the vim search command. To enter a search string, press the forward slash (/) key. The cursor goes to the message line, and vim displays a forward slash. Enter the text you want to find, and press the Enter key. The vim editor responds with one of three actions:

- If the word appears after the current cursor location, it jumps to the first location where the text appears.
- If the word doesn't appear after the current cursor location, it wraps around the end of the file to the first location in the file where the text appears (and indicates this with a message).
- It produces an error message stating that the text was not found in the file.

To continue searching for the same word, press the forward slash character and then press the Enter key, or you can use the n key, for *next*.

The substitute command allows you to quickly replace (substitute) one word for another in the text. To get to the substitute command you must be in command line mode. The format for the substitute command is:

```
:s/old/new/
```

The vim editor jumps to the first occurrence of the text `old` and replaces it with the text `new`. There are a few modifications you can make to the substitute command to substitute more than one occurrence of the text:

- `:s/old/new/g` to replace all occurrences of `old` in a line
- `:n,ms/old/new/g` to replace all occurrences of `old` between line numbers `n` and `m`
- `:%s/old/new/g` to replace all occurrences of `old` in the entire file
- `:%s/old/new/gc` to replace all occurrences of `old` in the entire file, but prompt for each occurrence

As you can see, for a command line text editor, vim contains quite a few advanced features. Because every Linux distribution includes it, it's a good idea to at least know the basics of the vim editor so that you can always edit scripts, no matter where you are or what you have available.

The emacs Editor

The emacs editor is an extremely popular editor that appeared before even Unix was around. Developers liked it so much they ported it to the Unix environment, and now it's been ported to the Linux environment. The emacs editor started out life as a console editor, much like `vi`, but has made the migration to the graphical world.

The emacs editor still provides the original console mode editor, but now it also has the ability to use a graphical X Windows window to allow editing text in a graphical environment. Typically, when you start the emacs editor from a command line, the editor will determine if you

have an available X Window session and start in graphical mode. If you don't, it will start in console mode.

This section describes both the console mode and graphical mode emacs editors so that you'll know how to use either one if you want (or need) to.

Using emacs on the Console

The console mode version of emacs is another editor that uses lots of key commands to perform editing functions. The emacs editor uses key combinations involving the Control key (the Ctrl key on the PC keyboard) and the Meta key. In most PC terminal emulator packages, the Meta key is mapped to the PC's Alt key. The official emacs documents abbreviate the Ctrl key as C- and the Meta key as M-. Thus, if you enter a Ctrl-x key combination, the document shows C-x. This chapter will do the same so as not to confuse you.

The Basics of emacs

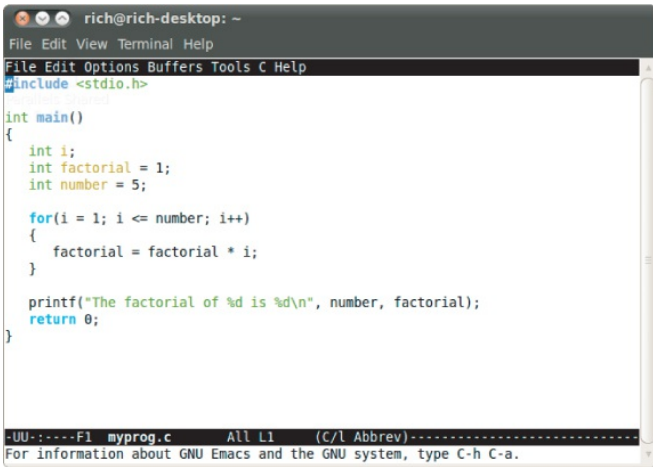
To edit a file using emacs, from the command line, enter:

```
$ emacs myprog.c
```

The emacs console mode window appears with a short introduction and help screen. Don't be alarmed; as soon as you press a key, emacs loads the file into the active buffer and displays the text, as shown in [Figure 9.2](#).

[Figure 9.2](#) Editing a file using the emacs editor in console

mode



The screenshot shows a terminal window titled "rich@rich-desktop: ~". The window contains the GNU Emacs editor in console mode. The menu bar at the top includes "File", "Edit", "View", "Terminal", and "Help". Below the menu bar, the editor's command line shows "File Edit Options Buffers Tools C Help". The main editing area displays a C program for calculating the factorial of a number. The status bar at the bottom shows "-UU-:----F1 myprog.c All L1 (C/l Abbrev)-----" and a message: "For information about GNU Emacs and the GNU system, type C-h C-a."

```
rich@rich-desktop: ~
File Edit View Terminal Help
File Edit Options Buffers Tools C Help
#include <stdio.h>
int main()
{
    int i;
    int factorial = 1;
    int number = 5;

    for(i = 1; i <= number; i++)
    {
        factorial = factorial * i;
    }

    printf("The factorial of %d is %d\n", number, factorial);
    return 0;
}
-UU-:----F1 myprog.c All L1 (C/l Abbrev)-----
For information about GNU Emacs and the GNU system, type C-h C-a.
```

You'll notice that the top of the console mode window shows a typical menu bar. Unfortunately, you won't be able to use the menu bar in console mode, only in graphical mode.

Note

If you have a graphical desktop but you prefer to use emacs in console mode instead of X Windows mode, use the `-nw` option on the command line.

Unlike the vim editor, where you have to move into and

out of insert mode to switch between entering commands and inserting text, the emacs editor has only one mode. If you type a printable character, emacs inserts it at the current cursor position. If you type a command, emacs executes the command.

To move the cursor around the buffer area, you can use the arrow keys and the PageUp and PageDown keys, assuming that emacs detected your terminal emulator correctly. If not, there are commands for moving the cursor around:

- `C-p` to move up one line (the previous line in the text).
- `C-b` to move left (back) one character.
- `C-f` to move right (forward) one character.
- `C-n` to move down one line (the next line in the text).

There are also commands for making longer jumps with the cursor within the text:

- `M-f` moves right (forward) to the next word.
- `M-b` moves left (backward) to the previous word.
- `C-a` moves to the beginning of the current line.
- `C-e` moves to the end of the current line.
- `M-a` moves to the beginning of the current sentence.
- `M-e` moves to the end of the current sentence.
- `M-v` moves back one screen of data.
- `C-v` moves forward one screen of data.
- `M-<` to move the first line of the text.
- `M->` to move to the last line of the text.

There are several commands you should know for saving

the editor buffer back into the file, and exiting emacs:

- `C-x C-s` to save the current buffer contents to the file.
- `C-z` to exit emacs but keep it running in your session so that you can come back to it.
- `C-x C-c` to exit emacs and stop the program.

You'll notice that two of these features require two key commands. The `C-x` command is called the *extend command*. This provides yet another whole set of commands to work with.

Editing Data

The emacs editor is pretty robust about inserting and deleting text in the buffer. To insert text, just move the cursor to the location where you want to insert the text and start typing. To delete text, emacs uses the Backspace key to delete the character before the current cursor position and the Delete key to delete the character at the current cursor location.

The emacs editor also has commands for killing text. The difference between deleting text and killing text is that when you kill text, emacs places it in a temporary area where you can retrieve it (see the “Copying and Pasting” section). Deleted text is gone forever.

There are a few commands for killing text in the buffer:

- `M-Backspace` to kill the word before the current cursor position
- `M-d` to kill the word after the current cursor position

- `C-k` to kill from the current cursor position to the end of the line
- `M-k` to kill from the current cursor position to the end of the sentence

The emacs editor also includes a fancy way of mass-killing text. Just move the cursor to the start of the area you want to kill and press either the `C-@` or `C-Spacebar` keys. Then move the cursor to the end of the area you want to kill and press the `C-w` command keys. All of the text between the two locations is killed.

If you happen to make a mistake when killing text, the `C-u` command will undo the kill command, and return the data the state it was in before you killed it.

Copying and Pasting

You've seen how to cut data from the emacs buffer area; now it's time to see how to paste it somewhere else. Unfortunately, if you use the vim editor, this process may confuse you when you use the emacs editor.

In an unfortunate coincidence, pasting data in emacs is called *yanking*. In the vim editor, copying is called yanking, which is what makes this a difficult thing to remember if you happen to use both editors.

After you kill data using one of the kill commands, move the cursor to the location where you want to paste the data, and use the `C-y` command. This yanks the text out of the temporary area and pastes it at the current cursor position. The `C-y` command yanks the text from the last kill command.

If you've performed multiple kill commands, you can cycle through them using the `M-y` command.

To copy text, just yank it back into the same location you killed it from and then move to the new location and use the `C-y` command again. You can yank text back as many times as you desire.

Searching and Replacing

Searching for text in the emacs editor is done by using the `C-s` and `C-r` commands. The `C-s` command performs a forward search in the buffer area from the current cursor position to the end of the buffer, whereas the `C-r` command performs a backward search in the buffer area from the current cursor position to the start of the buffer.

When you enter either the `C-s` or `C-r` command, a prompt appears in the bottom line, querying you for the text to search. There are two types of searches that emacs can perform.

In an *incremental* search, the emacs editor performs the text search in real-time mode as you type the word. When you type the first letter, it highlights all of the occurrences of that letter in the buffer. When you type the second letter, it highlights all of the occurrences of the two-letter combination in the text, and so on until you complete the text you're searching for.

In a *non-incremental* search, press the Enter key after the `C-s` or `C-r` commands. This locks the search query into

the bottom line area and allows you to type the search text in full before searching.

To replace an existing text string with a new text string, you have to use the `M-x` command. This command requires a text command, along with parameters.

The text command is `replace-string`. After typing the command, press the Enter key, and emacs will query you for the existing text string. After entering that, press the Enter key again, and emacs will query you for the new replacement text string.

Using Buffers in emacs

The emacs editor allows you to edit multiple files at the same time by having multiple buffer areas. You can load files into a buffer and switch between buffers while editing.

To load a new file into a buffer while you're in emacs, use the `C-x C-f` key combination. This is the emacs Find a File mode. It takes you to the bottom line in the window and allows you to enter the name of the file you want to start to edit. If you don't know the name or location of the file, just press the Enter key. This brings up a file browser in the edit window, as shown in [Figure 9.3](#).

[Figure 9.3](#) The emacs Find a File mode browser


```

/home/rich:
total used in directory 544 available 56310608
drwxr-xr-x 37 rich rich 4096 2010-10-07 19:38 .
drwxr-xr-x 4 root root 4096 2010-09-23 19:01 ..
drwx----- 2 rich rich 4096 2010-09-22 12:10 .aptitude
-rw----- 1 rich rich 2831 2010-09-27 20:32 .bash_history
-rw-r--r-- 1 rich rich 220 2010-05-10 18:36 .bash_logout
-rw-r--r-- 1 rich rich 3103 2010-05-10 18:36 .bashrc
drwx----- 5 rich rich 4096 2010-10-07 19:27 .cache
drwx----- 3 rich rich 4096 2010-09-11 18:34 .compiz
drwxr-xr-x 10 rich rich 4096 2010-09-11 18:32 .config
drwx----- 3 rich rich 4096 2010-05-10 18:47 .dbus
drwxr-xr-x 2 rich rich 4096 2010-09-22 12:10 .debtags
drwxr-xr-x 2 rich rich 4096 2010-09-11 18:31 Desktop
-rw-r--r-- 1 rich rich 41 2010-10-07 19:27 .dmrc
drwxr-xr-x 2 rich rich 4096 2010-05-10 18:47 Documents
drwxr-xr-x 2 rich rich 4096 2010-09-22 13:46 Downloads
drwxr-xr-x 3 rich rich 4096 2010-10-07 19:38 .emacs.d
-rw----- 1 rich rich 16 2010-05-10 18:47 .esd_auth
-rw-r--r-- 1 rich rich 179 2010-05-10 18:36 examples.desktop
drwxr-xr-x 2 rich rich 4096 2010-06-28 20:09 .fontconfig
-UUU:%%-F1 ~ Top L5 (Dired by name)-----

```

From here, you can browse to the file you want to edit. To traverse up a directory level, go to the double dot entry, and press the Enter key. To traverse down a directory, go to the directory entry and press the Enter key. When you've found the file you want to edit, just press the Enter key, and emacs will load it into a new buffer area.

You can list the active buffer areas by pressing the `C-x C-b` extended command combination. The emacs editor splits the editor window and displays a list of buffers in the bottom window. There are always two buffers that emacs provides in addition to your main editing buffer:

- A scratch area called `*scratch*`

- A message area called *Messages*

The scratch area allows you to enter LISP programming commands as well as enter notes to yourself. The message area shows messages generated by emacs while operating. If any errors occur while using emacs, they will appear in the message area.

There are two ways to switch to a different buffer area in the window:

- `C-x o` to switch to the buffer listing window. Use the arrow keys to move to the buffer area you want and press the Enter key.
- `C-x b` to type in the name of the buffer area you want to switch to.

When you select the option to switch to the buffer listing window, emacs will open the buffer area in the new window area. The emacs editor allows you to have multiple windows open in a single session. The following section discusses how to manage multiple windows in emacs.

Using Windows in Console Mode emacs

The console mode emacs editor was developed many years before the idea of graphical windows appeared. However, it was advanced for its time, in that it could support multiple editing windows within the main emacs window.

You can split the emacs editing window into multiple

windows by using one of two commands:

- `C-x 2` splits the window horizontally into two windows.
- `C-x 3` splits the window vertically into two windows.

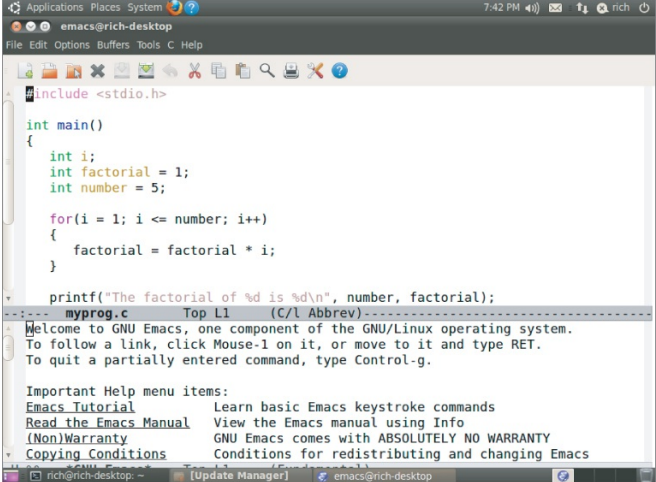
To move from one window to another, use the `C-x o` command. You'll notice that when you create a new window, emacs uses the buffer area from the original window in the new window. Once you move into the new window, you can use the `C-x C-f` command to load a new file, or one of the commands to switch to a different buffer area in the new window.

To close a window, move to it and use the `C-x 0` (that's a zero) command. If you want to close all of the windows except the one you're in, use the `C-x 1` (that's a numerical one) command.

Using emacs in X Windows

If you use emacs from an X Windows environment (such as the KDE or GNOME desktops), it will start in graphical mode, as shown in [Figure 9.4](#).

[Figure 9.4](#) The emacs graphical window



If you've already used emacs in console mode, you should be fairly familiar with the X Windows mode. All of the key commands are available as menu bar items. The emacs menu bar contains the following items:

- **File:** Allows you to open files in the window, create new windows, close windows, save buffers, and print buffers.
- **Edit:** Allows you to cut and copy selected text to the clipboard, paste clipboard data to the current cursor position, search for text, and replace text.
- **Options:** Provides settings for many more emacs

features, such as highlighting, word wrap, cursor type, and setting fonts.

- **Buffers:** Lists the current buffers available and allows you to easily switch between buffer areas.
- **Tools:** Provides access to the advanced features in emacs, such as the command line interface access, spell checking, comparing text between files (called diff), sending an e-mail message, calendar, and the calculator.
- **Help:** Provides the emacs manual online for access to help on specific emacs functions.

In addition to the normal graphical emacs menu bar items, there is often a separate item specific to the file type in the editor buffer. [Figure 9.4](#) shows opening a C program, so emacs provided a C menu item, allowing advanced settings for highlighting C syntax, and compiling, running, and debugging the code from a command prompt.

The graphical emacs window is an example of an older console application making the migration to the graphical world. Now that many Linux distributions provide graphical desktops (even on servers that don't need them), graphical editors are becoming more commonplace. Both of the popular Linux desktop environments (KDE and GNOME) have also provided graphical text editors specifically for their environments, which are covered in the rest of this chapter.

The KDE Family of Editors

If you're using a Linux distribution that uses the KDE desktop (see Chapter 1), there are a couple of options for you when it comes to text editors. The KDE project officially supports two different text editors:

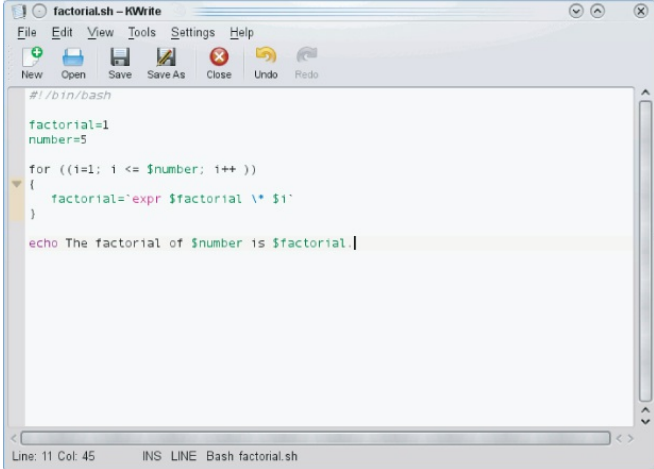
- **KWrite:** A single-screen text-editing package
- **Kate:** A full-featured, multi-window text-editing package

Both of these editors are graphical text editors that contain many advanced features. The Kate editor provides more advanced features, plus extra niceties not often found in standard text editors. This section describes each of the editors and shows some of the features that you can use to help with your shell script editing.

The KWrite Editor

The basic editor for the KDE environment is KWrite. It provides simple word-processing-style text editing, along with support for code syntax highlighting and editing. The default KWrite editing window is shown in [Figure 9.5](#).

Figure 9.5 The default KWrite window editing a shell script program



You can't tell from [Figure 9.5](#), but the KWrite editor recognizes several types of programming languages and uses color coding to distinguish constants, functions, and comments. Also, notice that the `for` loop has an icon that links the opening and closing braces. This is called a *folding marker*. By clicking the icon, you can collapse the function into a single line. This is a great feature when working through large applications.

The KWrite editing window provides full cut and paste capabilities, using the mouse and the arrow keys. As in a word processor, you can highlight and cut (or copy) text anywhere in the buffer area and paste it at any other place.

To edit a file using KWrite, you can either select KWrite from the KDE menu system on your desktop (some Linux distributions even create a Panel icon for it) or start it from the command line prompt:

```
$ kwrite factorial.sh
```

The `kwrite` command has several command line parameters you can use to customize how it starts:

- `--stdin` causes KWrite to read data from the standard input device instead of a file.
- `--encoding` specifies a character encoding type to use for the file.
- `--line` specifies a line number in the file to start at in the editor window.
- `--column` specifies a column number in the file to start at in the editor window.

The KWrite editor provides both a menu bar and a toolbar at the top of the edit window, allowing you to select features and change configuration settings of the KWrite editor.

The menu bar contains the following items:

- **File** to load, save, print, and export text from files.
- **Edit** to manipulate text in the buffer area.
- **View** to manage how the text appears in the editor window.
- **Bookmarks** for handling pointers to return to specific locations in the text (this option may need to be enabled in the configurations).

- **Tools** contains specialized features to manipulate the text.
- **Settings** for configuring the way the editor handles text.
- **Help** for getting information about the editor and commands.

The Edit menu bar item provides commands for all of your text-editing needs. Instead of having to remember cryptic key commands (which by the way, KWrite also supports), you can just select items in the Edit menu bar, as shown in [Table 9.2](#).

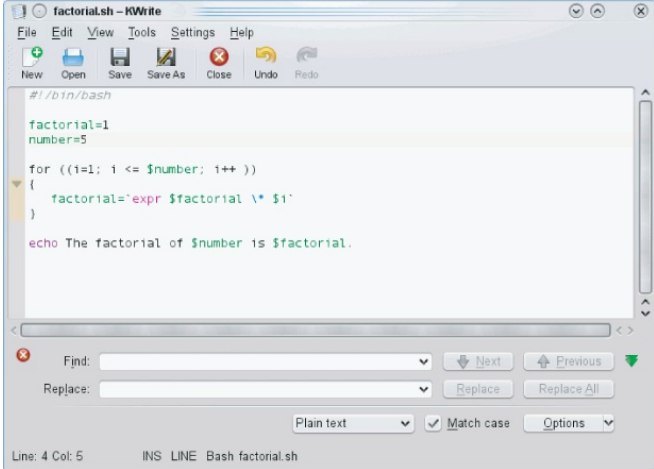
[Table 9.2](#) The KWrite Edit Menu Items

Item	Description
Undo	Reverses the last action or operation
Redo	Reverses the last undo action
Cut	Deletes the selected text and places it in the clipboard
Copy	Copies the selected text to the clipboard
Copy as HTML	Copies the selected text to the clipboard as HTML code
Paste	Inserts the current contents of the clipboard at the current cursor position
Select All	Selects all text in the editor
Deselect	Deselects any text that is currently selected
Overwrite Mode	Toggles insert mode to overwrite mode, replacing text with new typed text instead of just inserting the new text
Find	Produces the Find Text dialog box, which allows you to customize a text search

Find Next	Repeats the last find operation forward in the buffer area
Find Previous	Repeats the last find operation backwards in the buffer area
Replace	Produces the Replace With dialog box, which allows you to customize a text search and replace
Find Selected	Finds the next occurrence of the selected text
Find Selected Backwards	Finds the previous occurrence of the selected text
Go to Line	Produces the Goto dialog box, which allows you to enter a line number. The cursor moves to the specified line

The Find feature has two modes—a normal mode, which can perform simple text searches, and a power search and replace mode, where you can do advanced searching and replacing if necessary. You toggle between the two modes using the green arrow in the Find section, as shown in [Figure 9.6](#).

[Figure 9.6](#) The KWrite Find section



The Find power mode allows you to search not only with words, but with a regular expression (discussed in Chapter 19) for the search. There are a few other options you can use to customize the search as well, indicating, for example, whether or not to perform a case-sensitive search, or to look only for whole words instead of finding the text within words.

The Tools menu bar item provides several handy features for working with the text in the buffer area. [Table 9.3](#) describes the tools available in KWrite.

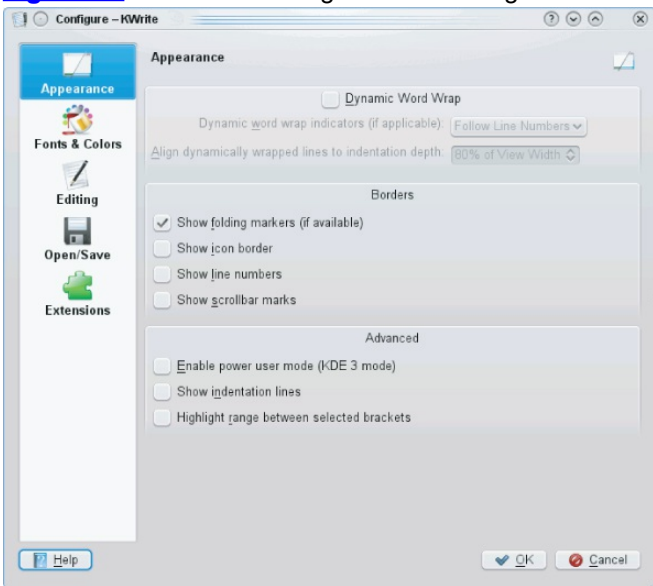
[Table 9.3](#) The KWrite Tools

Tool	Description
Read Only Mode	Locks the text so that no changes can be made while in the editor
Encoding	Sets the character set encoding used by the text
Spelling	Starts the spell-check program at the start of the text
Spelling (from cursor)	Starts the spell-check program from the current cursor position
Spellcheck Selection	Starts the spell-check program only on the selected section of text
Indent	Increases the paragraph indentation by one
Unindent	Decreases the paragraph indentation by one
Clean Indentation	Returns all paragraph indentation to the original settings
Align	Forces the current line or the selected lines to return to the default indentation settings
Uppercase	Sets the selected text, or the character at the current cursor position, to upper case
Lowercase	Sets the selected text, or the character at the current cursor position, to lower case
Capitalize	Capitalizes the first letter of the selected text or the word at the current cursor position
Join Lines	Combines the selected lines, or the line at the current cursor position, and the next line into one line
Word Wrap Document	Enable word wrapping in the text. If a line extends past the editor window edge, the line continues on the next line

There are lots of tools for a simple text editor!

The Settings menu includes the Configure Editor dialog box, shown in [Figure 9.7](#).

Figure 9.7 The KWrite Configure Editor dialog box



The Configuration dialog box uses icons on the left side for you to select the feature in KWrite to configure. When you select an icon, the right side of the dialog box shows the configuration settings for the feature.

The Appearance feature allows you to set several features that control how the text appears in the text editor window. You can enable word wrap, line numbers (great for programmers), and the folder markers from here. With the

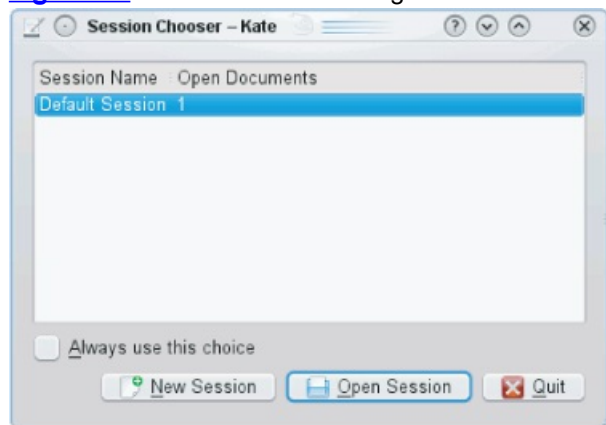
Fonts & Colors feature, you can customize the complete color scheme for the editor, determining what colors to make each category of text in the program code.

The Kate Editor

The Kate editor is the flagship editor for the KDE Project. It uses the same text editor as the KWrite application (so most of those features are the same), but it incorporates lots of other features into a single package.

The first thing you'll notice when you start the Kate editor is that the editor doesn't start! Instead, you get a dialog box, as shown in [Figure 9.8](#).

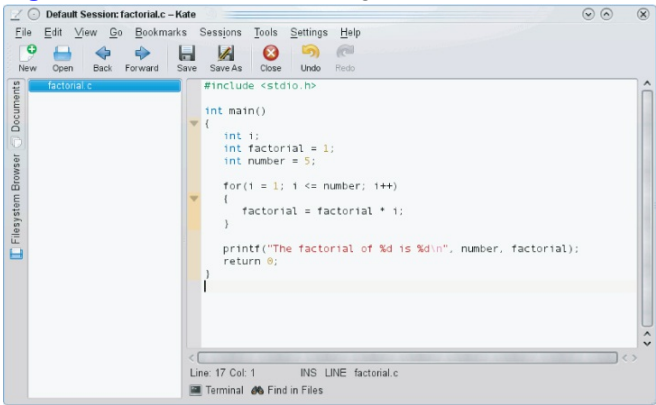
Figure 9.8 The Kate session dialog box



The Kate editor handles files in sessions. You can have multiple files open in a session, and you can have multiple sessions saved. When you start Kate, it provides you with the choice of which session to return to. When you close your Kate session, it remembers the documents you had open and displays them the next time you start Kate. This allows you to easily manage files from multiple projects by using separate workspaces for each project.

After selecting a session, you'll see the main Kate editor window, shown in [Figure 9.9](#).

Figure 9.9 The main Kate editing window

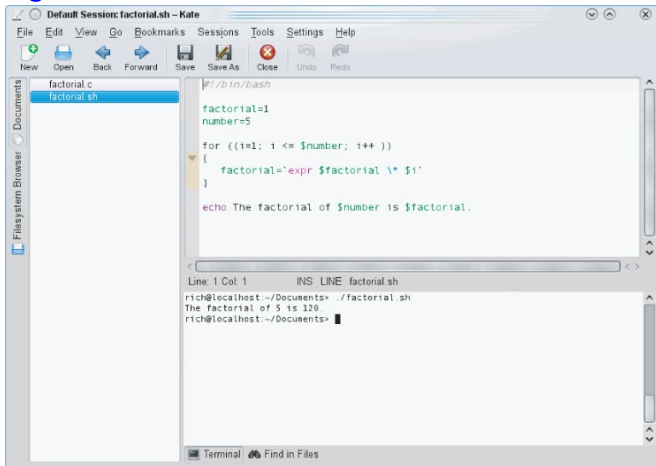


The left side frame shows the documents currently open in the session. You can switch between documents just by

clicking the document name. To edit a new file, click the Filesystem Browser tab on the left side. The left frame is now a full graphical filesystem browser, allowing you to graphically browse to locate your files.

A great feature of the Kate editor is the built-in terminal window, shown in [Figure 9.10](#).

Figure 9.10 The Kate built-in terminal window



The terminal tab at the bottom of the text editor window starts the built-in terminal emulator in Kate (using the KDE Konsole terminal emulator). This feature horizontally splits the current editing window, creating a new window with Konsole running in it. You can now enter command line

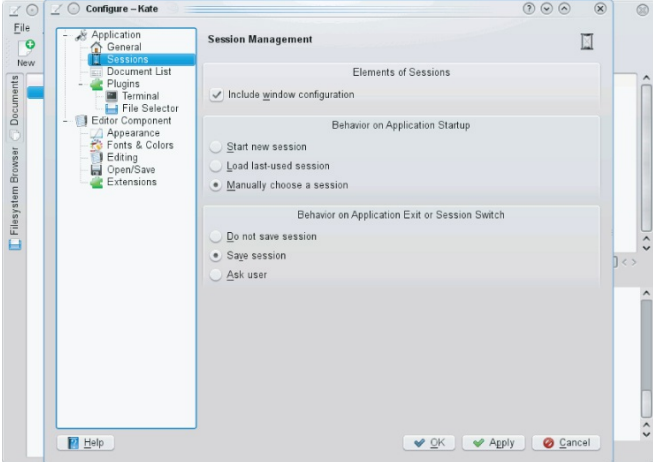
commands, start programs, or check on system settings without having to leave the editor! To close the terminal window, just type **exit** at the command prompt.

As you can tell from the terminal feature, Kate also supports multiple windows. The Window menu bar item provides options to:

- Create a new Kate window using the current session
- Split the current window vertically to create a new window
- Split the current window horizontally to create a new window
- Close the current window

To set the configuration settings in Kate, select the Configure Kate item under the Settings menu bar item. The Configuration dialog box, shown in [Figure 9.11](#), appears.

[Figure 9.11](#) The Kate configuration settings dialog box



You'll notice that the Editor settings area is exactly the same as for KWrite. This is because the two editors share the same text editor engine. The Application settings area allows you to configure settings for the Kate items, such as controlling sessions (shown in [Figure 9.11](#)), the documents list, and the filesystem browser. Kate also supports external plug-in applications, which can be activated here.

The GNOME Editor

If you're working on a Linux system using the GNOME desktop environment, there's a graphical text editor that you

can use as well. The gedit text editor is a basic text editor, with a few advanced features thrown in just for fun. This section walks you through the features of gedit and demonstrates how to use it for your shell script programming.

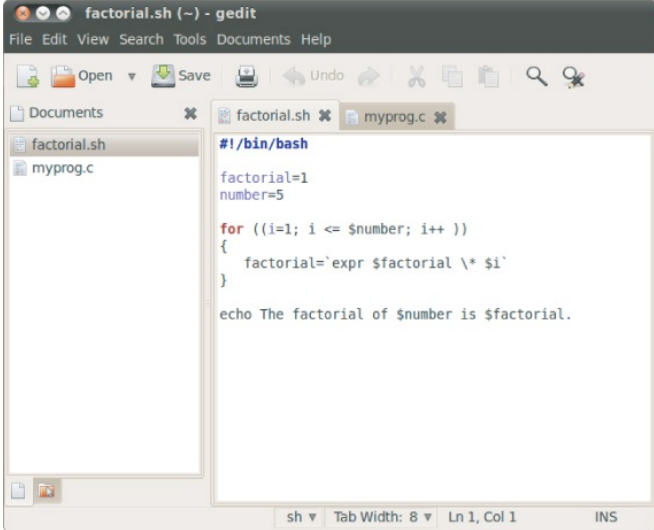
Starting gedit

Most GNOME desktop environments include gedit in the Accessories Panel menu item. If you can't find gedit there, you can start it from the command line prompt:

```
$ gedit factorial.sh myprog.c
```

When you start gedit with multiple files, it loads all of the files into separate buffers and displays each one as a tabbed window within the main editor window, as shown in [Figure 9.12](#).

[Figure 9.12](#) The gedit main editor window



The left frame in the gedit main editor window shows the documents you're currently editing. The right side shows the tabbed windows that contain the buffer text. If you hover your mouse pointer over each tab, a dialog box appears, showing the full pathname of the file, the MIME type, and the character set encoding it uses.

Basic gedit Features

In addition to the editor windows, gedit uses both a menu

bar and toolbar that allow you to set features and configure settings. The toolbar provides quick access to menu bar items. The menu bar items available are:

- **File** for handling new files, saving existing files, and printing files
- **Edit** to manipulate text in the active buffer area and set the editor preferences
- **View** to set the editor features to display in the window and to set the text highlighting mode
- **Search** to find and replace text in the active editor buffer area
- **Tools** to access plug-in tools installed in gedit
- **Documents** to manage files open in the buffer areas
- **Help** to access the full gedit manual

There shouldn't be anything too surprising here. The File menu provides the option Open Location, which allows you to open a file from the network using the standard Uniform Resource Identifier (URI) format popular in the World Wide Web world. This format identifies the protocol used to access the file (such as HTTP or FTP), the server where the file is located, and the complete path on the server to access the file.

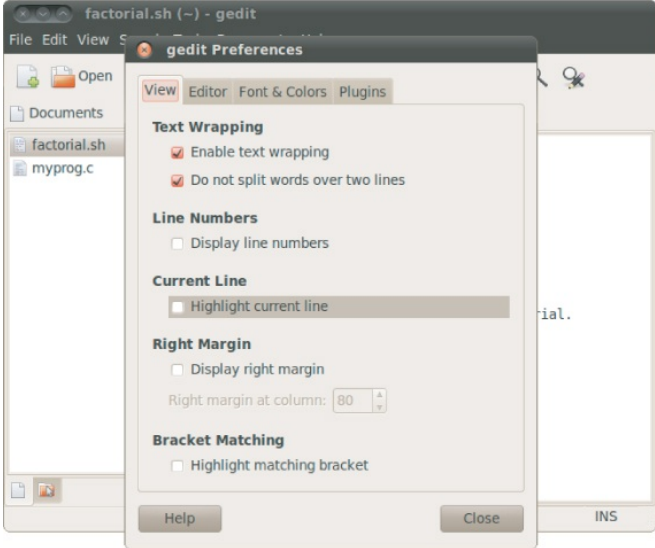
The Edit menu contains the standard cut, copy, and paste functions, along with a neat feature that allows you to easily enter the date and time in the text in several different formats. The Search menu provides a standard find function, which produces a dialog box where you can enter

the text to find, along with the capability to select how the find should work (matching case, matching the whole word, and the search direction). It also provides an incremental search feature, which works in real-time mode, finding text as you type the characters of the word.

Setting Preferences

The Edit menu contains a Preferences item, which produces the gedit Preferences dialog box, shown in [Figure 9.13](#).

[Figure 9.13](#) The gedit Preferences dialog box



This is where you can customize the operation of the gedit editor. The Preferences dialog box contains five tabbed areas for setting the features and behavior of the editor.

View

The View tab provides options for how gedit displays the text in the editor window:

- **Text Wrapping:** Determines how to handle long

lines of text in the editor. The Enabling text wrapping option wraps long lines to the next line of the editor. The Do Not Split Words Over Two Lines option prevents the auto-inserting of hyphens into long words, to prevent them being split between two lines.

- **Line Numbers:** Displays line numbers in the left margin in the editor window.
- **Current Line:** Highlights the line where the cursor is currently positioned, enabling you to easily find the cursor position.
- **Right Margin:** Enables the right margin and allows you to set how many columns should be in the editor window. The default value is 80 columns.
- **Bracket Matching:** When enabled, highlights bracket pairs in programming code, allowing you to easily match brackets in `if-then` statements, `for` and `while` loops, and other coding elements that use brackets.

The line-numbering and bracket-matching features provide an environment for programmers to troubleshoot code that's not often found in text editors.

Editor

The Editor tab provides options for how the gedit editor handles tabs and indentation, along with how files are saved:

- **Tab Stops:** Sets the number of spaces skipped

when you press the Tab key. The default value is eight. This feature also includes a checkbox that, when selected, inserts spaces instead of a tab skip.

- **Automatic Indentation:** When enabled, causes gedit to automatically indent lines in the text for paragraphs and code elements (such as `if-then` statements and loops).
- **File Saving:** Provides two features for saving files: whether or not to create a backup copy of the file when opened in the edit window, and whether or not to automatically save the file at a preselected interval.

The auto-save feature is a great way to ensure that your changes are saved on a regular basis to prevent catastrophes from crashes or power outages.

Font & Colors

The Font & Colors tab allows you to configure (not surprisingly) two items:

- **Font:** Allows you to select the default font of Monospace 10, or to select a customized font and font size from a dialog box.
- **Color Scheme:** Allows you to select the default color scheme used for text, background, selected text, and selection colors, or choose a custom color for each category.

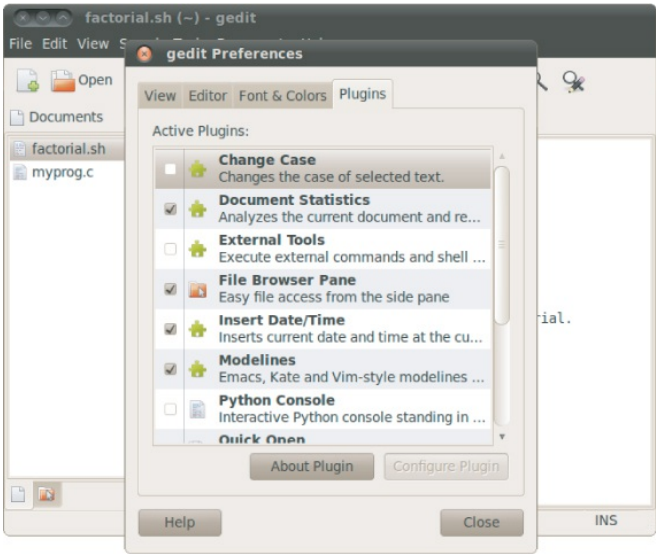
The default colors for gedit normally match the standard GNOME desktop theme selected for the desktop. These

colors will change to match the scheme you select for the desktop.

Plug-ins

The Plugins tab provides control over the plug-ins used in gedit. Plug-ins are separate programs that can interface with gedit to provide additional functionality. The Plugins tab is shown in [Figure 9.14](#).

Figure 9.14 The gedit Plugins Preferences tab



Several plug-ins are available for gedit, but not all of them are installed by default. [Table 9.4](#) describes the plug-ins that are currently available in gedit.

Table 9.4 The gedit Plug-ins

Plug-in	Description
Change Case	Changes the case of selected text
Document Statistics	Reports the number of words, lines, characters, and non-space characters
External Tools	Provides a shell environment in the editor to execute commands and scripts
File Browser Pane	Provides a simple file browser to make selecting files for editing easier
Insert Date/Time	Inserts the current date and time in several formats at the current cursor position
Modelines	Provides emacs-style message lines at the bottom of the editor window
Python Console	Provides an interactive console at the bottom of the editor window for entering commands using the Python programming language
Quick Open	Opens files directly in the gedit edit window
Snippets	Allows you to store often-used pieces of text for easy retrieval anywhere in the text
Sort	Quickly sorts the entire file or selected text
Spell Checker	Provides dictionary spellchecking for the text file
Tag List	Provides a list of commonly used strings you can easily enter into your text

Plug-ins that are enabled show a checkmark in the checkbox next to their name. Some plug-ins, such as the External Tools plug-in, also provide additional configuration features after you select them. It allows you to set a shortcut key to start the terminal, where gedit displays output, and the command to use to start the shell session.

Unfortunately, not all plug-ins are installed in the same place in the gedit menu bar. Some plug-ins appear in the Tools menu bar item (such as the Spell Checker and External Tools plug-ins), while others appear in the Edit menu bar item (such as the Change Case and Insert Date/Time plug-ins).

Summary

When it comes to creating shell scripts, you'll need some type of text editor. Several popular text editors are available for the Linux environment. The most popular editor in the Unix world, vi, has been ported to the Linux world as the vim editor. The vim editor provides simple text editing from the console, using a rudimentary full-screen graphical mode. The vim editor provides many advanced editor features, such as text searching and replacement.

Another popular Unix editor, emacs, has also made its way to the Linux world. The Linux version of emacs has both console and an X Windows graphical mode, making it the bridge between the old world and the new. The emacs

editor provides multiple buffer areas, allowing you to edit multiple files simultaneously.

The KDE Project created two editors for use in the KDE desktop. The KWrite editor is a simple editor that provides the basic text-editing features, along with a few advanced features such as syntax highlighting for programming code, line numbering, and code folding. The Kate editor provides more advanced features for programmers. One great feature in Kate is a built-in terminal window. You can open a command line interface session directly in the Kate editor without having to open a separate terminal emulator window. The Kate editor also allows you to open multiple files, providing different windows for each opened file.

The GNOME Project also provides a simple text editor for programmers. The gedit editor is a basic text editor that provides some advanced features such as code syntax highlighting and line numbering, but it was designed to be a bare-bones editor. To spruce up the gedit editor, developers created plug-ins, which expand the features available in gedit. Current plug-ins include a spell-checker, a terminal emulator, and a file browser.

This wraps up the background chapters on working with the command line in Linux. The next part of the book dives into the shell-scripting world. The next chapter starts off by showing you how to create a shell script file and how to run it on your Linux system. It will also show you the basics of shell scripts, allowing you to create simple programs by stringing multiple commands together into a script you can

run.

Part II

Shell Scripting Basics

In This Part

Chapter 10: Basic Script Building

Chapter 11: Using Structured Commands

Chapter 12: More Structured Commands

Chapter 13: Handling User Input

Chapter 14: Presenting Data

Chapter 15: Script Control

Chapter 10

Basic Script Building

In This Chapter

- Basic script building
- Using multiple commands
- Creating a script file

Now that we've covered the basics of the Linux system and the command line, it's time to start coding. This chapter discusses the basics of writing shell scripts. You'll need to know these basic concepts before you can start writing your own shell script masterpieces.

Using Multiple Commands

So far you've seen how to use the command line interface (CLI) prompt of the shell to enter commands and view the command results. The key to shell scripts is the ability to enter multiple commands and process the results from each command, even possibly passing the results of one command to another. The shell allows you to chain

commands together into a single step.

If you want to run two commands together, you can enter them on the same prompt line, separated with a semicolon:

```
$ date ; who
Mon Feb 21 15:36:09 EST 2011
Christine tty2          2011-02-21 15:26
Samantha tty3          2011-02-21 15:26
Timothy tty1           2011-02-21 15:26
user tty7              2011-02-19 14:03 (:0)
user pts/0             2011-02-21 15:21 (:0.0)

$
```

Congratulations, you just wrote a shell script! This simple script uses just two bash shell commands. The `date` command runs first, displaying the current date and time, followed by the output of the `who` command, showing who is currently logged on to the system. Using this technique, you can string together as many commands as you wish, up to the maximum command line character count of 255 characters.

While using this technique is fine for small scripts, it has a major drawback in that you have to enter the entire command at the command prompt every time you want to run it. Instead of having to manually enter the commands onto a command line, you can combine the commands into a simple text file. When you need to run the commands, just simply run the text file.

Creating a Script File

To place shell commands in a text file, first you'll need to

use a text editor (see Chapter 9) to create a file, then enter the commands into the file.

When creating a shell script file, you must specify the shell you are using in the first line of the file. The format for this is:

```
#!/bin/bash
```

In a normal shell script line, the pound sign (#) is used as a comment line. A comment line in a shell script isn't processed by the shell. However, the first line of a shell script file is a special case, and the pound sign followed by the exclamation point tells the shell what shell to run the script under (yes, you can be using a bash shell and run your script using another shell).

After indicating the shell, commands are entered onto each line of the file, followed by a carriage return. As mentioned, comments can be added by using the pound sign. An example looks like this:

```
#!/bin/bash
# This script displays the date and who's logged on
date
who
```

And that's all there is to it. You can use the semicolon and put both commands on the same line if you want to, but in a shell script, you can list commands on separate lines. The shell will process commands in the order in which they appear in the file.

Also notice that another line was included that starts with the pound symbol and adds a comment. Lines that start with the pound symbol (other than the first #! line) aren't

interpreted by the shell. This is a great way to leave comments for yourself about what's happening in the script, so when you come back to it two years later you can easily remember what you did.

Save this script in a file called `test1`, and you are almost ready. There are still a couple of things to do before you can run your new shell script file.

If you try running the file now, you'll be somewhat disappointed to see this:

```
$ test1
bash: test1: command not found
$
```

The first hurdle to jump is getting the bash shell to find your script file. If you remember from Chapter 5, the shell uses an environment variable called `PATH` to find commands. A quick look at the `PATH` environment variable demonstrates our problem:

```
$ echo $PATH
/usr/kerberos/sbin:/usr/kerberos/bin:/usr/local/bin:/usr/bin
:/bin:/usr/local/sbin:/usr/sbin:/sbin:/home/user/bin $
```

The `PATH` environment variable is set to look for commands only in a handful of directories. To get the shell to find the `test1` script, we need to do one of two things:

- Add the directory where our shell script file is located to the `PATH` environment variable.
- Use an absolute or relative filepath to reference our shell script file in the prompt.

Tip

Some Linux distributions add the `$HOME/bin` directory to the `PATH` environment variable. This creates a place in every user's `HOME` directory to place files where the shell can find them to execute.

For this example, we'll use the second method to tell the shell exactly where the script file is located. Remember that to reference a file in the current directory, you can use the single dot operator in the shell:

```
$ ./test1
bash: ./test1: Permission denied
$
```

Now the shell found the shell script file just fine, but there's another problem. The shell indicated that you don't have permission to execute the file. A quick look at the file permissions should show what's going on here:

```
$ ls -l test1
-rw-r--r--  1 user      user           73 Sep 24 19:56 test1
$
```

When the new `test1` file was created, the `umask` value determined the default permission settings for the new file. Because the `umask` variable is set to 022 (see Chapter 6), the system created the file with only read/write permissions for the file's owner.

The next step is to give the file owner permission to execute the file, using the `chmod` command (see Chapter 6):

```
$ chmod u+x test1
$ ./test1
Mon Feb 21 15:38:19 EST 2011
Christine tty2          2011-02-21 15:26
Samantha tty3           2011-02-21 15:26
Timothy tty1            2011-02-21 15:26
user tty7               2011-02-19 14:03 (:0)
```

Success! Now all of the pieces are in the right places to execute the new shell script file.

Displaying Messages

Most shell commands produce their own output, which is displayed on the console monitor where the script is running. Many times, however, you will want to add your own text messages to help the script user know what is happening within the script. You can do this with the `echo` command. The `echo` command can display a simple text string if you add the string following the command:

```
$ echo This is a test
This is a test
$
```

Notice that by default you don't need to use quotes to delineate the string you're displaying. However, sometimes this can get tricky if you are using quotes within your string:

```
$ echo Let's see if this'll work
Lets see if thisll work
$
```

The `echo` command uses either double or single quotes to delineate text strings. If you use them within your string, you need to use one type of quote within the text and the other type to delineate the string:

```
$ echo "This is a test to see if you're paying attention"
This is a test to see if you're paying attention
$ echo 'Rich says "scripting is easy".'
Rich says "scripting is easy".
$
```

Now all of the quotation marks appear properly in the output.

You can add `echo` statements anywhere in your shell scripts where you need to display additional information:

```
$ cat test1
#!/bin/bash
# This script displays the date and who's logged on
echo The time and date are:
date
echo "Let's see who's logged into the system:"
who
$
```

When you run this script, it produces the following output:

```
$ ./test1
The time and date are:
Mon Feb 21 15:41:13 EST 2011
Let's see who's logged into the system:
Christine tty2          2011-02-21 15:26
Samantha tty3          2011-02-21 15:26
Timothy tty1           2011-02-21 15:26
user tty7              2011-02-19 14:03 (:0)
user pts/0             2011-02-21 15:21 (:0.0)
$
```

That's nice, but what if you want to `echo` a text string on the same line as a command output? You can use the `-n` parameter for the `echo` statement to do that. Just change the first `echo` statement line to this:

```
echo -n "The time and date are: "
```

You'll need to use quotes around the string to ensure that there's a space at the end of the echoed string. The command output begins exactly where the string output stops. The output will now look like this:

```
$ ./test1
The time and date are: Mon Feb 21 15:42:23 EST 2011
Let's see who's logged into the system:
Christine tty2          2011-02-21 15:26
Samantha tty3          2011-02-21 15:26
```

```
Timothy tty1 2011-02-21 15:26
user tty7 2011-02-19 14:03 (:0)
user pts/0 2011-02-21 15:21 (:0.0)
$
```

Perfect! The `echo` command is a crucial piece of shell scripts that interact with users. You'll find yourself using it in many situations, especially when you want to display the values of script variables. Let's look at that next.

Using Variables

Just running individual commands from the shell script is useful, but this has its limitations. Often you'll want to incorporate other data in your shell commands to process information. You can do this by using *variables*. Variables allow you to temporarily store information within the shell script for use with other commands in the script. This section shows how to use variables in your shell scripts.

Environment Variables

You've already seen one type of Linux variable in action. Chapter 5 described the environment variables available in the Linux system. You can access these values from your shell scripts as well.

The shell maintains environment variables that track specific system information, such as the name of the system, the name of the user logged in to the system, the user's system ID (called UID), the default home directory of

the user, and the search path used by the shell to find programs. You can display a complete list of active environment variables available by using the `set` command:

```
$ set
BASH=/bin/bash
...
HOME=/home/Samantha
HOSTNAME=localhost.localdomain
HOSTTYPE=i386
IFS=$' \t\n'
IMSETTINGS_INTEGRATE_DESKTOP=yes
IMSETTINGS_MODULE=none
LANG=en_US.utf8
LESSOPEN='|/usr/bin/lesspipe.sh %s'
LINES=24
LOGNAME=Samantha
...
```

You can tap into these environment variables from within your scripts by using the environment variable's name preceded by a dollar sign. This is demonstrated in the following script:

```
$ cat test2
#!/bin/bash
# display user information from the system.
echo "User info for userid: $USER"
echo UID: $UID
echo HOME: $HOME
$
```

The `$USER`, `$UID`, and `$HOME` environment variables are used to display the pertinent information about the logged-in user. The output should look something like this:

```
$ chmod u+x test2
$ ./test2
User info for userid: Samantha
UID: 1001
HOME: /home/Samantha
$ $
```

Notice that the environment variables in the `echo`

commands are replaced by their current values when the script is run. Also notice that we were able to place the `$USER` system variable within the double quotation marks in the first string, and the shell script was still able to figure out what we meant. There is a drawback to using this method, however. Look at what happens in this example:

```
$ echo "The cost of the item is $15"  
The cost of the item is 5
```

That is obviously not what was intended. Whenever the script sees a dollar sign within quotes, it assumes you're referencing a variable. In this example the script attempted to display the variable `$1` (which was not defined), and then the number `5`. To display an actual dollar sign, you must precede it with a backslash character:

```
$ echo "The cost of the item is \$15"  
The cost of the item is $15
```

That's better. The backslash allowed the shell script to interpret the dollar sign as an actual dollar sign, and not a variable. The next section shows how to create your own variables in your scripts.

Note

You may also see variables referenced using the format `${variable}`. The extra braces around the variable name are often used to help identify the variable name from the dollar sign.

User Variables

In addition to the environment variables, a shell script allows you to set and use your own variables within the script. Setting variables allows you to temporarily store data and use it throughout the script, making the shell script more like a real computer program.

User variables can be any text string of up to 20 letters, digits, or an underscore character. User variables are case sensitive, so the variable `Var1` is different from the variable `var1`. This little rule often gets novice script programmers in trouble.

Values are assigned to user variables using an equal sign. No spaces can appear between the variable, the equal sign, and the value (another trouble spot for novices). Here are a few examples of assigning values to user variables:

```
var1=10
var2=-57
var3=testing
var4="still more testing"
```

The shell script automatically determines the data type used for the variable value. Variables defined within the shell script maintain their values throughout the life of the shell script but are deleted when the shell script completes.

Just like system variables, user variables can be referenced using the dollar sign:

```
$ cat test3
#!/bin/bash
# testing variables
days=10
guest="Katie"
echo "$guest checked in $days days ago"
days=5
guest="Jessica"
```

```
echo "$guest checked in $days days ago"  
$
```

Running the script produces the following output:

```
$ chmod u+x test3  
$ ./test3  
Katie checked in 10 days ago  
Jessica checked in 5 days ago  
$
```

Each time the variable is referenced, it produces the value currently assigned to it. It's important to remember that when referencing a variable value you use the dollar sign, but when referencing the variable to assign a value to it, you do not use the dollar sign. Here's an example of what I mean:

```
$ cat test4  
#!/bin/bash  
# assigning a variable value to another variable  
  
value1=10  
value2=$value1  
echo The resulting value is $value2  
$
```

When you use the value of the `value1` variable in the assignment statement, you must still use the dollar sign. This code produces the following output:

```
$ chmod u+x test4  
$ ./test4  
The resulting value is 10  
$
```

If you forget the dollar sign, and make the `value2` assignment line look like:

```
value2=value1
```

you get the following output:

```
$ ./test4  
The resulting value is value1
```

Without the dollar sign the shell interprets the variable name as a normal text string, which is most likely not what you wanted.

The Backtick

One of the most useful features of shell scripts is the lowly back quote character, usually called the backtick (`) in the Linux world. Be careful—this is not the normal single quotation mark character you are used to using for strings. Because it is not used very often outside of shell scripts, you may not even know where to find it on your keyboard. You should become familiar with it, because it's a crucial component of many shell scripts. Hint: On a U.S. keyboard, it is usually on the same key as the tilde symbol (~).

The backtick allows you to assign the output of a shell command to a variable. While this doesn't seem like much, it is a major building block in script programming.

You must surround the entire command line command with backtick characters:

```
testing=`date`
```

The shell runs the command within the backticks and assigns the output to the variable testing. Here's an example of creating a variable using the output from a normal shell command:

```
$ cat test5
#!/bin/bash
# using the backtick character
testing=`date`
```

```
echo "The date and time are: " $testing
$
```

The variable `testing` receives the output from the `date` command, and it is used in the `echo` statement to display it. Running the shell script produces the following output:

```
$ chmod u+x test5
$ ./test5
The date and time are:  Mon Jan 31 20:23:25 EDT 2011
$
```

That's not all that exciting in this example (you could just as easily just put the command in the `echo` statement), but once you capture the command output in a variable, you can do anything with it.

Here's a popular example of how the backtick is used to capture the current date and use it to create a unique filename in a script:

```
#!/bin/bash
# copy the /usr/bin directory listing to a log file
today=`date +%y%m%d`
ls /usr/bin -al > log.$today
```

The `today` variable is assigned the output of a formatted date command. This is a common technique used to extract date information for log filenames. The `+%y%m%d` format instructs the date command to display the date as a two-digit year, month, and day:

```
$ date +%y%m%d
110131
$
```

The script assigns the value to a variable, which is then used as part of a filename. The file itself contains the redirected output (discussed later in the “Redirecting Input and Output” section) of a directory listing. After running the

script, you should see a new file in your directory:

```
-rw-r--r--  1 user  user    769 Jan 31 10:15 log.110131
```

The log file appears in the directory using the value of the `$today` variable as part of the filename. The contents of the log file are the directory listing from the `/usr/bin` directory. If the script is run the next day, the log filename will be `log.110201`, thus creating a new file for the new day.

Redirecting Input and Output

There are times when you'd like to save the output from a command instead of just having it displayed on the monitor. The bash shell provides a few different operators that allow you to *redirect* the output of a command to an alternative location (such as a file). Redirection can be used for input as well as output, redirecting a file to a command for input. This section describes what you need to do to use redirection in your shell scripts.

Output Redirection

The most basic type of redirection is sending output from a command to a file. The bash shell uses the greater-than symbol (`>`) for this:

```
command > outputfile
```

Anything that would appear on the monitor from the command instead is stored in the output file specified:

```
$ date > test6
```

```
ls -l test6
-rw-r--r-- 1 user user 29 Feb 10 17:56 test6
$ cat test6
Thu Feb 10 17:56:58 EDT 2011
$
```

The redirect operator created the file `test6` (using the default `umask` settings) and redirected the output from the `date` command to the `test6` file. If the output file already exists, the redirect operator overwrites the existing file with the new file data:

```
$ who > test6
$ cat test6
user pts/0 Feb 10 17:55
$
```

Now the contents of the `test6` file contain the output from the `who` command.

Sometimes, instead of overwriting the file's contents, you may need to append output from a command to an existing file, for example if you're creating a log file to document an action on the system. In this situation, you can use the double greater-than symbol (`>>`) to append data:

```
$ date >> test6
$ cat test6
user pts/0 Feb 10 17:55
Thu Feb 10 18:02:14 EDT 2011
$
```

The `test6` file still contains the original data from the `who` command processed earlier—plus now it contains the new output from the `date` command.

Input Redirection

Input redirection is the opposite of output redirection.

Instead of taking the output of a command and redirecting it to a file, input redirection takes the content of a file and redirects it to a command.

The input redirection symbol is the less-than symbol (<):

```
command < inputfile
```

The easy way to remember this is that the command is always listed first in the command line, and the redirection symbol “points” to the way the data is flowing. The less-than symbol indicates that the data is flowing from the input file to the command.

Here's an example of using input redirection with the `wc` command:

```
$ wc < test6
      2      11      60
$
```

The `wc` command provides a count of text in the data. By default, it produces three values:

- The number of lines in the text
- The number of words in the text
- The number of bytes in the text

By redirecting a text file to the `wc` command, you can get a quick count of the lines, words, and bytes in the file. The example shows that there are 2 lines, 11 words, and 60 bytes in the `test6` file.

There's another method of input redirection, called *inline input redirection*. This method allows you to specify the data for input redirection on the command line instead of in a file. This may seem somewhat odd at first, but there are a

few applications for this process (such as those shown in the “Performing Math” section later).

The inline input redirection symbol is the double less-than symbol (<<). Besides this symbol, you must specify a text marker that delineates the beginning and end of the data used for input. You can use any string value for the text marker, but it must be the same at the beginning of the data and the end of the data:

```
command << marker
data
marker
```

When using inline input redirection on the command line, the shell will prompt for data using the secondary prompt, defined in the `PS2` environment variable (see Chapter 5). Here's how this looks when you use it:

```
$ wc << EOF
> test string 1
> test string 2
> test string 3
> EOF
      3          9         42
$
```

The secondary prompt continues to prompt for more data until you enter the string value for the text marker. The `wc` command performs the line, word, and byte counts of the data supplied by the inline input redirection.

Pipes

There are times when you need to send the output of one command to the input of another command. This is

possible using redirection, but somewhat clunky:

```
$ rpm -qa > rpm.list
$ sort < rpm.list
abrt-1.1.14-1.fc14.i686
abrt-addon-ccpp-1.1.14-1.fc14.i686
abrt-addon-kerneloops-1.1.14-1.fc14.i686
abrt-addon-python-1.1.14-1.fc14.i686
abrt-desktop-1.1.14-1.fc14.i686
abrt-gui-1.1.14-1.fc14.i686
abrt-libs-1.1.14-1.fc14.i686
abrt-plugin-bugzilla-1.1.14-1.fc14.i686
abrt-plugin-logger-1.1.14-1.fc14.i686
abrt-plugin-runapp-1.1.14-1.fc14.i686
acl-2.2.49-8.fc14.i686

...
```

The `rpm` command manages the software packages installed on systems using the Red Hat Package Management system (RPM), such as the Fedora system as shown. When used with the `-qa` parameters, it produces a list of the existing packages installed, but not necessarily in any specific order. If you're looking for a specific package, or group of packages, it can be difficult to find it using the output of the `rpm` command.

Using the standard output redirection, the output was redirected from the `rpm` command to a file, called `rpm.list`. After the command finished, the `rpm.list` file contained a list of all the installed software packages on my system. Next, input redirection was used to send the contents of the `rpm.list` file to the `sort` command to sort the package names alphabetically.

That was useful, but again, a somewhat clunky way of producing the information. Instead of redirecting the output of a command to a file, you can redirect the output to

another command. This process is called piping.

Like the backtick (`), the symbol for piping is not used often outside of shell scripting. The symbol is two vertical lines, one above the other. However, the `pipe` symbol often looks like a single vertical line in print (`|`). On a U.S. keyboard, it is usually on the same key as the backslash (`\`). The `pipe` is put between the commands to redirect the output from one to the other:

```
command1 | command2
```

Don't think of piping as running two commands back to back. The Linux system actually runs both commands at the same time, linking them together internally in the system. As the first command produces output, it's sent immediately to the second command. No intermediate files or buffer areas are used to transfer the data.

Now, using piping you can easily `pipe` the output of the `rpm` command directly to the `sort` command to produce your results:

```
$ rpm -qa | sort
abrt-1.1.14-1.fc14.i686
abrt-addon-ccpp-1.1.14-1.fc14.i686
abrt-addon-kerneloops-1.1.14-1.fc14.i686
abrt-addon-python-1.1.14-1.fc14.i686
abrt-desktop-1.1.14-1.fc14.i686
abrt-gui-1.1.14-1.fc14.i686
abrt-libs-1.1.14-1.fc14.i686
abrt-plugin-bugzilla-1.1.14-1.fc14.i686
abrt-plugin-logger-1.1.14-1.fc14.i686
abrt-plugin-runapp-1.1.14-1.fc14.i686
acl-2.2.49-8.fc14.i686
...
```

Unless you're a (very) quick reader, you probably couldn't keep up with the output generated by this command.

Because the piping feature operates in real time, as soon as the `rpm` command produces data, the `sort` command gets busy sorting it. By the time the `rpm` command finishes outputting data, the `sort` command already has the data sorted and starts displaying it on the monitor.

There's no limit to the number of `pipe`s you can use in a command. You can continue piping the output of commands to other commands to refine your operation.

In this case, because the output of the `sort` command zooms by so quickly, you can use one of the text paging commands (such as `less` or `more`) to force the output to stop at every screen of data:

```
$ rpm -qa | sort | more
```

This command sequence runs the `rpm` command, `pipe`s the output to the `sort` command, and then `pipe`s that output to the `more` command to display the data, stopping after every screen of information. This now lets you pause and read what's on the display before continuing, as shown in [Figure 10.1](#).

[Figure 10.1](#) Using piping to send data to the `more` command

```
File Edit View Search Terminal Help
abrt-1.1.14-1.fc14.i686
abrt-addon-ccpp-1.1.14-1.fc14.i686
abrt-addon-kerneloops-1.1.14-1.fc14.i686
abrt-addon-python-1.1.14-1.fc14.i686
abrt-desktop-1.1.14-1.fc14.i686
abrt-gui-1.1.14-1.fc14.i686
abrt-libs-1.1.14-1.fc14.i686
abrt-plugin-bugzilla-1.1.14-1.fc14.i686
abrt-plugin-logger-1.1.14-1.fc14.i686
abrt-plugin-runapp-1.1.14-1.fc14.i686
acl-2.2.49-8.fc14.i686
alsa-firmware-1.0.23-1.fc14.noarch
alsa-lib-1.0.23-2.fc14.i686
alsa-plugins-pulseaudio-1.0.22-1.fc13.i686
alsa-tools-firmware-1.0.23-1.fc14.i686
alsa-utils-1.0.23-3.fc14.i686
anaconda-14.22-1.fc14.i686
anaconda-yum-plugins-1.0-5.fc12.noarch
anthy-9100h-15.fc14.i686
apr-1.3.9-3.fc13.i686
apr-util-1.3.10-1.fc14.i686
apr-util-ldap-1.3.10-1.fc14.i686
ar9170-firmware-2009.05.28-2.fc13.noarch
--More--
```

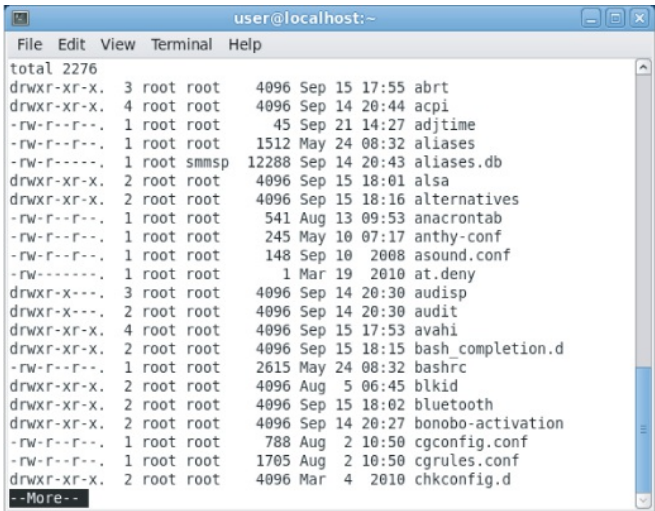
To get even fancier, you can use redirection along with piping to save your output to a file:

```
$ rpm -qa | sort > rpm.list
$ more rpm.list
abrt-1.1.14-1.fc14.i686
abrt-addon-ccpp-1.1.14-1.fc14.i686
abrt-addon-kerneloops-1.1.14-1.fc14.i686
abrt-addon-python-1.1.14-1.fc14.i686
abrt-desktop-1.1.14-1.fc14.i686
abrt-gui-1.1.14-1.fc14.i686
abrt-libs-1.1.14-1.fc14.i686
abrt-plugin-bugzilla-1.1.14-1.fc14.i686
abrt-plugin-logger-1.1.14-1.fc14.i686
abrt-plugin-runapp-1.1.14-1.fc14.i686
acl-2.2.49-8.fc14.i686
...
```

As expected, the data in the `rpm.list` file is now sorted!

By far one of the most popular uses of piping is piping the results of commands that produce long output to the `more` command. This is especially common with the `ls` command, as shown in [Figure 10.2](#).

Figure 10.2 Using the `more` command with the `ls` command



The `ls -l` command produces a long listing of all the files in the directory. For directories with lots of files, this can be quite a listing. By piping the output to the `more` command, you force the output to stop at the end of every screen of data.

Performing Math

Another feature crucial to any programming language is the ability to manipulate numbers. Unfortunately, for shell scripts this process is a bit awkward. There are two different ways to perform mathematical operations in your shell scripts.

The expr Command

Originally, the Bourne shell provided a special command that was used for processing mathematical equations. The `expr` command allowed the processing of equations from the command line, but it is extremely clunky:

```
$ expr 1 + 5
6
```

The `expr` command recognizes a few different mathematical and string operators, shown in [Table 10.1](#).

Table 10.1 The `expr` Command Operators

Operator	Description
<code>ARG1 ARG2</code>	Return <code>ARG1</code> if neither argument is null or zero; otherwise, return <code>ARG2</code> .
<code>ARG1 & ARG2</code>	Return <code>ARG1</code> if neither argument is null or zero; otherwise, return 0.
<code>ARG1 < ARG2</code>	Return 1 if <code>ARG1</code> is less than <code>ARG2</code> ; otherwise, return 0.
<code>ARG1 <= ARG2</code>	Return 1 if <code>ARG1</code> is less than or equal to <code>ARG2</code> ; otherwise, return 0.
<code>ARG1 = ARG2</code>	Return 1 if <code>ARG1</code> is equal to <code>ARG2</code> ; otherwise, return 0.
<code>ARG1 != ARG2</code>	Return 1 if <code>ARG1</code> is not equal to <code>ARG2</code> ; otherwise, return 0.
	Return 1 if <code>ARG1</code> is greater than or equal to <code>ARG2</code> ; otherwise,

<code>ARG1 >= ARG2</code>	return 0.
<code>ARG1 > ARG2</code>	Return 1 if <code>ARG1</code> is greater than <code>ARG2</code> ; otherwise, return 0.
<code>ARG1 + ARG2</code>	Return the arithmetic sum of <code>ARG1</code> and <code>ARG2</code> .
<code>ARG1 - ARG2</code>	Return the arithmetic difference of <code>ARG1</code> and <code>ARG2</code> .
<code>ARG1 * ARG2</code>	Return the arithmetic product of <code>ARG1</code> and <code>ARG2</code> .
<code>ARG1 / ARG2</code>	Return the arithmetic quotient of <code>ARG1</code> divided by <code>ARG2</code> .
<code>ARG1 % ARG2</code>	Return the arithmetic remainder of <code>ARG1</code> divided by <code>ARG2</code> .
<code>STRING : REGEXP</code>	Return the pattern match if <code>REGEXP</code> matches a pattern in <code>STRING</code> .
<code>match STRING REGEXP</code>	Return the pattern match if <code>REGEXP</code> matches a pattern in <code>STRING</code> .
<code>substr STRING POS LENGTH</code>	Return the substring <code>LENGTH</code> characters in length, starting at position <code>POS</code> (starting at 1).
<code>index STRING CHARS</code>	Return position in <code>STRING</code> where <code>CHARS</code> is found; otherwise, return 0.
<code>length STRING</code>	Return the numeric length of the string <code>STRING</code> .
<code>+ TOKEN</code>	Interpret <code>TOKEN</code> as a string, even if it's a keyword.
<code>(EXPRESSION)</code>	Return the value of <code>EXPRESSION</code> .

While the standard operators work fine in the `expr` command, the problem occurs when using them from a script or the command line. Many of the `expr` command operators have other meanings in the shell (such as the asterisk). Using them in the `expr` command produces odd results:

```
$ expr 5 * 2
expr: syntax error
$
```

To solve this problem, you need to use the shell escape character (the backslash) to identify any characters that

may be misinterpreted by the shell before being passed to the `expr` command:

```
$ expr 5 \* 2
10
$
```

Now that's really starting to get ugly! Using the `expr` command in a shell script is equally cumbersome:

```
$ cat test6
#!/bin/bash
# An example of using the expr command
var1=10
var2=20
var3=`expr $var2 / $var1`
echo The result is $var3
```

To assign the result of a mathematical equation to a variable, you have to use the backtick character to extract the output from the `expr` command:

```
$ chmod u+x test6
$ ./test6
The result is 2
$
```

Fortunately, the bash shell has an improvement for processing mathematical operators as you shall see in the next section.

Using Brackets

The bash shell includes the `expr` command to stay compatible with the Bourne shell; however, it also provides a much easier way of performing mathematical equations. In bash, when assigning a mathematical value to a variable, you can enclose the mathematical equation using a dollar sign and square brackets (`${ operation }`):

```
$ var1=$((1 + 5))
$ echo $var1
6
$ var2 = [$var1 * 2]
$ echo $var2
12
$
```

Using brackets makes shell math much easier than with the `expr` command. This same technique also works in shell scripts:

```
$ cat test7
#!/bin/bash
var1=100
var2=50
var3=45
var4=$((var1 * (var2 - var3))
echo The final result is $var4
$
```

Running this script produces the output:

```
$ chmod u+x test7
$ ./test7
The final result is 500
$
```

Also, notice that when using the square brackets method for calculating equations you don't need to worry about the multiplication symbol, or any other characters, being misinterpreted by the shell. The shell knows that it's not a wildcard character because it is within the square brackets.

There's one major limitation to performing math in the bash shell script. Take a look at this example:

```
$ cat test8
#!/bin/bash
var1=100
var2=45
var3=$((var1 / var2))
echo The final result is $var3
$
```

Now run it and see what happens:

```
$ chmod u+x test8
$ ./test8
The final result is 2
$
```

The bash shell mathematical operators support only integer arithmetic. This is a huge limitation if you're trying to do any sort of real-world mathematical calculations.

Note

The z shell (zsh) provides full floating-point arithmetic operations. If you require floating-point calculations in your shell scripts, you might consider checking out the z shell (discussed in Chapter 22).

A Floating-Point Solution

There are several solutions for overcoming the bash integer limitation. The most popular solution uses the built-in bash calculator, called `bc`.

The Basics of bc

The bash calculator is actually a programming language that allows you to enter floating-point expressions at a command line and then interprets the expressions, calculates them, and returns the result. The bash calculator recognizes:

- Numbers (both integer and floating point)

- Variables (both simple variables and arrays)
- Comments (lines starting with a pound sign or the C language `/* */` pair)
- Expressions
- Programming statements (such as `if-then` statements)
- Functions

You can access the bash calculator from the shell prompt using the `bc` command:

```
$ bc
bc 1.06.95
Copyright 1991-1994, 1997, 1998, 2000, 2004, 2006 Free Software Foundation, Inc.
This is free software with ABSOLUTELY NO WARRANTY.
For details type 'warranty'.
12 * 5.4
64.8
3.156 * (3 + 5)
25.248
quit
$
```

The example starts out by entering the expression `12 * 5.4`. The bash calculator returns the answer. Each subsequent expression entered into the calculator is evaluated, and the result is displayed. To exit the bash calculator, you must enter `quit`.

The floating-point arithmetic is controlled by a built-in variable called `scale`. You must set this value to the desired number of decimal places you want in your answers or you won't get what you were looking for:

```
$ bc -q
3.44 / 5
0
scale=4
3.44 / 5
.6880
quit
```

The default value for the `scale` variable is zero. Before the `scale` value is set, the bash calculator provides the answer to zero decimal places. After you set the `scale` variable value to four, the bash calculator displays the answer to four decimal places. The `-q` command line parameter suppresses the lengthy welcome banner from the bash calculator.

In addition to normal numbers, the bash calculator also understands variables:

```
$ bc -q
var1=10
var1 * 4
40
var2 = var1 / 5
print var2
2
quit
$
```

Once a variable value is defined, you can use the variable throughout the bash calculator session. The `print` statement allows you to print variables and numbers.

Using bc in Scripts

Now you may be wondering how the bash calculator is going to help you with floating-point arithmetic in your shell scripts. Do you remember your friend the backtick character? Yes, you can use the backtick character to run a `bc` command and assign the output to a variable! The basic format to use is this:

```
variable='echo "options; expression" | bc'
```

The first portion, `options`, allows you to set variables. If you need to set more than one variable, separate them using the semicolon. The `expression` parameter defines the mathematical expression to evaluate using `bc`. Here's a quick example of doing this in a script:

```
$ cat test9
#!/bin/bash
var1='echo " scale=4; 3.44 / 5" | bc'
echo The answer is $var1
$
```

This example sets the `scale` variable to four decimal places and then specifies a specific calculation for the expression. Running this script produces the following output:

```
$ chmod u+x test9
$ ./test9
The answer is .6880
$
```

Now that's fancy! You aren't limited to just using numbers for the expression value. You can also use variables defined in the shell script:

```
$ cat test10
#!/bin/bash
var1=100
var2=45
var3='echo "scale=4; $var1 / $var2" | bc'
echo The answer for this is $var3
$
```

The script defines two variables, which are used within the expression sent to the `bc` command. Remember to use the dollar sign to signify the value for the variables and not the variables themselves. The output of this script is as follows:

```
$ ./test10
```

The answer for this is 2.2222

\$

And of course, once a value is assigned to a variable, that variable can be used in yet another calculation:

```
$ cat test11
#!/bin/bash
var1=20
var2=3.14159
var3='echo "scale=4; $var1 * $var1" | bc'
var4='echo "scale=4; $var3 * $var2" | bc'
echo The final result is $var4
$
```

This method works fine for short calculations, but sometimes you need to get more involved with your numbers. If you have more than just a couple of calculations, it gets confusing trying to list multiple expressions on the same command line.

There's a solution to this problem. The `bc` command recognizes input redirection, allowing you to redirect a file to the `bc` command for processing. However, this also can get confusing, as you'd need to store your expressions in a file.

The best method is to use inline input redirection, which allows you to redirect data directly from the command line. In the shell script, you assign the output to a variable:

```
variable='bc << EOF
options
statements
expressions
EOF
'
```

The `EOF` text string indicates the beginning and end of the inline redirection data. Remember that the backtick characters are still needed to assign the output of the `bc`

command to the variable.

Now you can place all of the individual bash calculator elements on separate lines in the script file. Here's an example of using this technique in a script:

```
$ cat test12
#!/bin/bash

var1=10.46
var2=43.67
var3=33.2
var4=71

var5='bc << EOF
scale = 4
a1 = ( $var1 * $var2)
b1 = ($var3 * $var4)
a1 + b1
EOF
'

echo The final answer for this mess is $var5
$
```

Placing each option and expression on a separate line in your script makes things cleaner and easier to read and follow. The `EOF` string indicates the start and end of the data to redirect to the `bc` command. Of course, you need to use the backtick characters to indicate the command to assign to the variable.

You'll also notice in this example that you can assign variables within the bash calculator. It's important to remember that any variables created within the bash calculator are valid only within the bash calculator and can't be used in the shell script.

Exiting the Script

So far in our sample scripts, we terminated things pretty abruptly. When we were done with our last command, we just ended the script. There's a more elegant way of completing things available to us.

Every command that runs in the shell uses an `exit status` to indicate to the shell that it's done processing. The exit status is an integer value between 0 and 255 that's passed by the command to the shell when the command finishes running. You can capture this value and use it in your scripts.

Checking the exit Status

Linux provides the `$?` special variable that holds the exit status value from the last command that executed. You must view or use the `$?` variable immediately after the command you want to check. It changes values to the exit status of the last command executed by the shell:

```
$ date
Sat Jan 15 10:01:30 EDT 2011
$ echo $?
0
$
```

By convention, the exit status of a command that successfully completes is zero. If a command completes with an error, then a positive integer value is placed in the exit status:

```
$ asdfg
-bash: asdfg: command not found
```

```
echo $?  
127  
$
```

The invalid command returns an exit status of 127. There's not much of a standard convention to Linux error exit status codes. However, there are a few guidelines you can use, as shown in [Table 10.2](#).

Table 10.2 Linux Exit Status Codes

Code	Description
0	Successful completion of the command
1	General unknown error
2	Misuse of shell command
126	The command can't execute
127	Command not found
128	Invalid exit argument
128+x	Fatal error with Linux signal x
130	Command terminated with Ctrl+C
255	Exit status out of range

An exit status value of 126 indicates that the user didn't have the proper permissions set to execute the command:

```
./myprog.c  
-bash: ./myprog.c: Permission denied  
echo $?  
126  
$
```

Another common error you'll encounter occurs if you supply an invalid parameter to a command:

```
date %t  
date: invalid date '%t'  
echo $?
```

```
$
```

This generates the general exit status code of one, indicating an unknown error occurred in the command.

The exit Command

By default, your shell script will exit with the exit status of the last command in your script:

```
$ ./test6
The result is 2
$ echo $?
0
$
```

You can change that to return your own exit status code. The `exit` command allows you to specify an exit status when your script ends:

```
$ cat test13
#!/bin/bash
# testing the exit status
var1=10
var2=30
var3=$(( $var1 + var2 ))
echo The answer is $var3
exit 5
$
```

When you check the exit status of the script, you'll get the value used as the parameter of the `exit` command:

```
$ chmod u+x test13
$ ./test13
The answer is 40
$ echo $?
5
$
```

You can also use variables in the `exit` command parameter:

```
$ cat test14
#!/bin/bash
# testing the exit status
var1=10
var2=30
var3=$(( $var1 + var2 ))
exit $var3
$
```

When you run this command, it produces the following exit status:

```
$ chmod u+x test14
$ ./test14
$ echo $?
40
$
```

You should be careful with this feature, however, as the exit status codes can only go up to 255. Watch what happens in this example:

```
$ cat test14b
#!/bin/bash
# testing the exit status
var1=10
var2=30
var3=$(( $var1 * var2 ))
echo The value is $var3
exit $var3
$
```

Now when you run it, you get the following:

```
$ ./test14b
The value is 300
$ echo $?
44
$
```

The exit status code is reduced to fit in the 0 to 255 range. The shell does this by using modulo arithmetic. The *modulo* of a value is the remainder after a division. The resulting number is the remainder of the specified number divided by 256. In the case of 300 (the result value), the

remainder is 44, which is what appears as the exit status code.

In Chapter 11, you'll see how you can use the `if-then` statement to check the error status returned by a command to see if the command was successful or not.

Summary

The bash shell script allows you to string commands together into a script. The most basic way to create a script is to separate multiple commands on the command line using a semicolon. The shell executes each command in order, displaying the output of each command on the monitor.

You can also create a shell script file, placing multiple commands in the file for the shell to execute in order. The shell script file must define the shell used to run the script. This is done in the first line of the script file, using the `#!` symbol, followed by the full path of the shell.

Within the shell script you can reference environment variable values by using a dollar sign in front of the variable. You can also define your own variables for use within the script, and assign values and even the output of a command by using the backtick character. The variable value can be used within the script by placing a dollar sign in front of the variable name.

The bash shell allows you to redirect both the input and

output of a command from the standard behavior. You can redirect the output of any command from the monitor display to a file by using the greater-than symbol, followed by the name of the file to capture the output. You can append output data to an existing file by using two greater-than symbols. The less-than symbol is used to redirect input to a command. You can redirect input from a file to a command.

The Linux `pipe` command (the broken bar symbol) allows you to redirect the output of a command directly to the input of another command. The Linux system runs both commands at the same time, sending the output of the first command to the input of the second command without using any redirect files.

The bash shell provides a couple of ways for you to perform mathematical operations in your shell scripts. The `expr` command is a simple way to perform integer math. In the bash shell, you can also perform basic math calculations by enclosing equations in square brackets, preceded by a dollar sign. To perform floating-point arithmetic, you need to utilize the `bc` calculator command, redirecting input from inline data and storing the output in a user variable.

Finally, the chapter discussed how to use the exit status in your shell script. Every command that runs in the shell produces an exit status. The exit status is an integer value between 0 and 255 that indicates if the command completed successfully, and if not, what the reason may

have been. An exit status of 0 indicates that the command completed successfully. You can use the `exit` command in your shell script to declare a specific exit status upon the completion of your script.

So far in your shell scripts, things have proceeded in an orderly fashion from one command to the next. In the next chapter, you'll see how you can use some logic flow control to alter which commands are executed within the script.

Chapter 11

Using Structured Commands

In This Chapter

- Working with the if-then statement
- The if-then-else statement
- Nesting ifs
- The test command
- Compound condition testing
- Advanced if-then features
- The code command
- Managing user accounts

In the shell scripts presented in Chapter 10, the shell processed each individual command in the shell script in the order it appeared. This works out fine for sequential operations, where you want all of the commands to process in the proper order. However, this isn't how all programs operate.

Many programs require some sort of logic flow control between the commands in the script. This means that the shell executes certain commands given one set of

circumstances, but it has the ability to execute other commands given a different set of circumstances. There is a whole class of commands that allows the script to skip over or loop through commands based on conditions of variable values or the result of other commands. These commands are generally referred to as *structured commands*.

The structured commands allow you to alter the flow of operation of a program, executing some commands under some conditions while skipping others under other conditions. There are quite a few structured commands available in the bash shell, so we'll look at them individually. In this chapter, we look at the `if-then` statement.

Working with the if-then Statement

The most basic type of structured command is the `if-then` statement. The `if-then` statement has the following format:

```
if command
then
    commands
fi
```

If you're using `if-then` statements in other programming languages, this format may be somewhat confusing. In other programming languages, the object after the `if` statement is an equation that is evaluated for a `TRUE` or `FALSE` value. That's not how the bash shell `if` statement works.

The bash shell `if` statement runs the command defined on the `if` line. If the exit status of the command (see Chapter 10) is zero (the command completed successfully), the commands listed under the `then` section are executed. If the exit status of the command is anything else, the `then` commands aren't executed, and the bash shell moves on to the next command in the script.

Here's a simple example to demonstrate this concept:

```
$ cat test1
#!/bin/bash
# testing the if statement
if date
then
    echo "it worked"
fi
```

This script uses the `date` command on the `if` line. If the command completes successfully, the `echo` statement should display the text string. When you run this script from the command line, you'll get the following results:

```
$ ./test1
Sat Jan 23 14:09:24 EDT 2011
it worked
$
```

The shell executed the `date` command listed on the `if` line. Since the exit status was zero, it also executed the `echo` statement listed in the `then` section.

Here's another example:

```
$ cat test2
#!/bin/bash
# testing a bad command
if asdfg
then
    echo "it did not work"
fi
echo "we are outside of the if statement"
```

```
$ ./test2
./test2: line 3: asdfg: command not found
we are outside of the if statement
$
```

In this example, a command was deliberately used that will not work in the `if` statement line. Because this is a bad command, it will produce an exit status that's non-zero, and the bash shell skips the `echo` statement in the `then` section. Also notice that the error message generated from running the command in the `if` statement still appears in the output of the script. There will be times when you won't want this to happen. Chapter 14 discusses how this can be avoided.

You are not limited to just one command in the `then` section. You can list commands just as in the rest of the shell script. The bash shell treats the commands as a block, executing all of them when the command in the `if` statement line returns a zero exit status or skipping all of them when the command returns a non-zero exit status:

```
$ cat test3
#!/bin/bash
# testing multiple commands in the then section
testuser=rich
if grep $testuser /etc/passwd
then
    echo The bash files for user $testuser are:
    ls -a /home/$testuser/.b*
fi
```

The `if` statement line uses the `grep` command to search the `/etc/passwd` file to see if a specific username is currently used on the system. If there's a user with that logon name, the script displays some text, and then lists the bash files in the user's `HOME` directory:

```
$ ./test3
```

```
rich:x:500:500:Rich Blum:/home/rich:/bin/bash
```

The files for user rich are:

```
/home/rich/.bash_history  /home/rich/.bash_profile
/home/rich/.bash_logout  /home/rich/.bashrc
$
```

However, if you set the `testuser` variable to a user that doesn't exist on the system, nothing happens:

```
$ ./test3
$
```

That's not all that exciting. It would be nice if we could display a little message saying that the username wasn't found on the system. Well, we can, using another feature of the `if-then` statement.

Note

You might see an alternative form of the `if-then` statement used in some scripts:

```
if command; then
  commands
fi
```

By putting a semicolon at the end of the command to evaluate, you can include the `then` statement on the same line, which looks more like how `if-then` statements are handled in some other programming languages.

The if-then-else Statement

In the `if-then` statement, you have only one option of whether or not a command is successful. If the command returns a non-zero exit status code, the bash shell just moves on to

the next command in the script. In this situation, it would be nice to be able to execute an alternate set of commands. That's exactly what the `if-then-else` statement is for.

The `if-then-else` statement provides another group of commands in the statement:

```
if command
then
    commands
else
    commands
fi
```

When the command in the `if` statement line returns with an exit status code of zero, the commands listed in the `then` section are executed, just as in a normal `if-then` statement. When the command in the `if` statement line returns a non-zero exit status code, the `bash` shell executes the commands in the `else` section.

Now you can modify the test script to look like this:

```
$ cat test4
#!/bin/bash
# testing the else section
testuser=badtest
if grep $testuser /etc/passwd
then
    echo The files for user $testuser are:
    ls -a /home/$testuser/.b*
else
    echo "The user name $testuser does not exist on this system"
fi
$
$ ./test4
The user name badtest does not exist on this system
$
```

That's more user-friendly. Just like the `then` section, the `else` section can contain multiple commands. The `fi` statement delineates the end of the `else` section.

Nesting ifs

Sometimes you must check for several situations in your script code. Instead of having to write separate `if-then` statements, you can use an alternative version of the `else` section, called `elif`.

The `elif` continues an `else` section with another `if-then` statement:

```
if command1
then
    commands
elif command2
then
    more commands
fi
```

The `elif` statement line provides another command to evaluate, similar to the original `if` statement line. If the exit status code from the `elif` command is zero, `bash` executes the commands in the second `then` statement section.

You can continue to string `elif` statements together, creating one huge `if-then-elif` conglomeration:

```
if command1
then
    command set 1
elif command2
then
    command set 2
elif command3
then
    command set 3
elif command4
then
    command set 4
fi
```

Each block of commands is executed depending on which command returns the zero exit status code. Remember, the bash shell will execute the `if` statements in order, and only the first one that returns a zero exit status will result in the `then` section being executed. Later on in “The case Command” section, you’ll see how to use the `case` command instead of having to nest lots of `if-then` statements.

The test Command

So far, all you’ve seen in the `if` statement line are normal shell commands. You might be wondering if the bash `if-then` statement has the ability to evaluate any condition other than the exit status code of a command.

The answer is no, it can’t. However, there’s a neat utility available in the bash shell that helps you evaluate other things, using the `if-then` statement.

The `test` command provides a way to test different conditions in an `if-then` statement. If the condition listed in the `test` command evaluates to true, the `test` command exits with a zero exit status code, making the `if-then` statement behave in much the same way that `if-then` statements work in other programming languages. If the condition is false, the `test` command exits with a1, which causes the `if-then` statement to fail.

The format of the test command is pretty simple:

```
test condition
```

The *condition* is a series of parameters and values that the `test` command evaluates. When used in an `if-then` statement, the `test` command looks like this:

```
if test condition
then
    commands
fi
```

The bash shell provides an alternative way of declaring the `test` command in an `if-then` statement:

```
if [ condition ]
then
    commands
fi
```

The square brackets define the condition that's used in the `test` command. Be careful; you *must* have a space after the first bracket and a space before the last bracket or you'll get an error message.

The `test` command can evaluate three classes of conditions:

- Numeric comparisons
- String comparisons
- File comparisons

The next sections describe how to use each of these classes of tests in your `if-then` statements.

Numeric Comparisons

The most common method for using the `test` command is to perform a comparison of two numeric values. [Table 11.1](#) shows the list of condition parameters used for testing two

values.

Table 11.1 The test Numeric Comparisons

Comparison	Description
<code>n1 -eq n2</code>	Check if <code>n1</code> is equal to <code>n2</code> .
<code>n1 -ge n2</code>	Check if <code>n1</code> is greater than or equal to <code>n2</code> .
<code>n1 -gt n2</code>	Check if <code>n1</code> is greater than <code>n2</code> .
<code>n1 -le n2</code>	Check if <code>n1</code> is less than or equal to <code>n2</code> .
<code>n1 -lt n2</code>	Check if <code>n1</code> is less than <code>n2</code> .
<code>n1 -ne n2</code>	Check if <code>n1</code> is not equal to <code>n2</code> .

The numeric test conditions can be used to evaluate both numbers and variables. Here's an example of doing that:

```
$ cat test5
#!/bin/bash
# using numeric test comparisons
vall=10
val2=11

if [ $vall -gt 5 ]
then
    echo "The test value $vall is greater than 5"
fi

if [ $vall -eq $val2 ]
then
    echo "The values are equal"
else
    echo "The values are different"
fi
```

The first test condition:

```
if [ $vall -gt 5 ]
```

tests if the value of the variable `val1` is greater than 5. The second test condition:

```
if [ $vall -eq $val2 ]
```

tests if the value of the variable *val1* is equal to the value of the variable *val2*. Run the script and watch the results:

```
$ ./test5
The test value 10 is greater than 5
The values are different
$
```

Both of the numeric test conditions evaluated as expected.

There is a limitation to the test numeric conditions, however. Try this script:

```
$ cat test6
#!/bin/bash
# testing floating point numbers
vall=' echo "scale=4; 10 / 3 " | bc\'
echo "The test value is $vall"
if [ $vall -gt 3 ]
then
    echo "The result is larger than 3"
fi
$
$ ./test6
The test value is 3.3333
./test6: line 5: [: 3.3333: integer expression expected
$
```

This example uses the bash calculator to produce a floating-point value, stored in the *val1* variable. Next, it uses the `test` command to evaluate the value. Something obviously went wrong here.

In Chapter 10, you learned how to trick the bash shell into handling floating-point values; there's still a problem in this script. The `test` command wasn't able to handle the floating-point value that was stored in the *val1* variable.

Remember that the only numbers the bash shell can handle are integers. When you utilize the bash calculator, you just fool the shell into storing a floating-point value in a

variable as a string value. This works perfectly fine if all you need to do is display the result, using an `echo` statement, but this doesn't work in numeric-oriented functions, such as our numeric test condition. The bottom line is that you're not able to use floating-point values in the `test` command.

String Comparisons

The `test` command also allows you to perform comparisons on string values. Performing comparisons on strings can get tricky, as you'll see. [Table 11.2](#) shows the comparison functions you can use to evaluate two string values.

Table 11.2 The test Command String Comparisons

Comparison	Description
<code>str1 = str2</code>	Check if <code>str1</code> is the same as string <code>str2</code> .
<code>str1 != str2</code>	Check if <code>str1</code> is not the same as <code>str2</code> .
<code>str1 < str2</code>	Check if <code>str1</code> is less than <code>str2</code> .
<code>str1 > str2</code>	Check if <code>str1</code> is greater than <code>str2</code> .
<code>-n str1</code>	Check if <code>str1</code> has a length greater than zero.
<code>-z str1</code>	Check if <code>str1</code> has a length of zero.

The following sections describe the different string comparisons available.

String Equality

The `equal` and `not equal` conditions are fairly self-explanatory with strings. It's pretty easy to know when two

string values are the same or not:

```
$cat test7
#!/bin/bash
# testing string equality
testuser=rich

if [ $USER = $testuser ]
then
    echo "Welcome $testuser"
fi
$
$ ./test7
Welcome rich
$
```

Also, using the not equals string comparison, will allow you to determine if two strings have the same value or not:

```
$ cat test8
#!/bin/bash
# testing string equality
testuser=baduser

if [ $USER != $testuser ]
then
    echo "This is not $testuser"
else
    echo "Welcome $testuser"
fi
$
$ ./test8
This is not baduser
$
```

The `test` comparison takes all punctuation and capitalization into account when comparing strings for equality.

String Order

Trying to determine if one string is less than or greater than another is where things start getting tricky. There are two problems that often plague shell programmers when trying

to use the greater-than or less-than features of the `test` command:

- The greater-than and less-than symbols must be escaped, or the shell will use them as redirection symbols, with the string values as filenames.
- The greater-than and less-than order is not the same as that used with the `sort` command.

The first item can result in a huge problem that often goes undetected when programming your scripts. Here's a typical example of what sometimes happens to novice shell script programmers:

```
$ cat badtest
#!/bin/bash
# mis-using string comparisons

vall=baseball
val2=hockey

if [ $vall > $val2 ]
then
    echo "$vall is greater than $val2"
else
    echo "$vall is less than $val2"
fi
$
$ ./badtest
baseball is greater than hockey
$ ls -l hockey
-rw-r--r--  1 rich      rich          0 Sep 30 19:08 hockey
$
```

By just using the greater-than symbol itself in the script, no errors are generated, but the results are wrong. The script interpreted the greater-than symbol as an output redirection. Thus, it created a file called `hockey`. Because the redirection completed successfully, the `test` command returns a zero exit status code, which the `if` statement

evaluates as though things completed successfully!

To fix this problem, you need to properly escape the greater-than symbol:

```
$ cat test9
#!/bin/bash
# mis-using string comparisons

vall=baseball
val2=hockey

if [ $vall \> $val2 ]
then
    echo "$vall is greater than $val2"
else
    echo "$vall is less than $val2"
fi
$
$ ./test9
baseball is less than hockey
$
```

Now that answer is more along the lines of what you would expect from the string comparison.

The second issue is a little more subtle, and you may not even run across it unless you are working with uppercase and lowercase letters. The `sort` command handles uppercase letters opposite to the way the `test` command considers them. Let's test this feature in a script:

```
$ cat test9b
#!/bin/bash
# testing string sort order

vall=Testing
val2=testing

if [ $vall \> $val2 ]
then
    echo "$vall is greater than $val2"
else
    echo "$vall is less than $val2"
fi
$
$ ./test9b
```

```
Testing is less than testing
$ sort testfile
testing
Testing
$
```

Capitalized letters are treated as less than lowercase letters in the `test` command. However, when you put the same strings in a file and use the `sort` command, the lowercase letters appear first. This is due to the ordering technique each command uses. The `test` command uses standard ASCII ordering, using each character's ASCII numeric value to determine the sort order. The `sort` command uses the sorting order defined for the system locale language settings. For the English language, the locale settings specify that lowercase letters appear before uppercase letters in sorted order.

Caution

Notice that the `test` command uses the standard mathematical comparison symbols for string comparisons and text codes for numerical comparisons. This is a subtle feature that many programmers manage to get reversed. If you use the mathematical comparison symbols for numeric values, the shell interprets them as string values and may not produce the correct results.

String Size

The `-n` and `-z` comparisons are handy when trying to evaluate if a variable contains data or not:

```
$ cat test10
```

```
#!/bin/bash
# testing string length
val1=testing
val2=''

if [ -n $val1 ]
then
    echo "The string '$val1' is not empty"
else
    echo "The string '$val1' is empty"
fi

if [ -z $val2 ]
then
    echo "The string '$val2' is empty"
else
    echo "The string '$val2' is not empty"
fi

if [ -z $val3 ]
then
    echo "The string '$val3' is empty"
else
    echo "The string '$val3' is not empty"
fi

$
$ ./test10
The string 'testing' is not empty
The string '' is empty
The string '' is empty
$
```

This example creates two string variables. The `val1` variable contains a string, and the `val2` variable is created as an empty string. The following comparisons are made as shown below:

```
if [ -n $val1 ]
```

determines if the *val1* variable is non-zero in length, which it is, so the `then` section is processed:

```
if [ -z $var2 ]
```

determines if the *val2* variable is zero in length, which it is, so the `then` section is processed:


```
if [ -z $val3 ]
```

determines if the `val3` variable is zero in length. This variable was never defined in the shell script, so it indicates that the string length is still zero, even though it wasn't defined.

Caution

Empty and uninitialized variables can have catastrophic effects on your shell script tests. If you're not sure of the contents of a variable, it's always best to test if the variable contains a value using `-n` or `-z` before using it in a numeric or string comparison.

File Comparisons

The last category of test comparisons is quite possibly the most powerful and most used comparisons in shell scripting. The `test` command allows you to test the status of files and directories on the Linux filesystem. [Table 11.3](#) lists these comparisons.

Table 11.3 The test Command File Comparisons

Comparison	Description
<code>-d file</code>	Check if <code>file</code> exists and is a directory.
<code>-e file</code>	Checks if <code>file</code> exists.
<code>-f file</code>	Checks if <code>file</code> exists and is a file.
<code>-r file</code>	Checks if <code>file</code> exists and is readable.
<code>-s file</code>	Checks if <code>file</code> exists and is not empty.

<code>-w file</code>	Checks if <code>file</code> exists and is writable.
<code>-x file</code>	Checks if <code>file</code> exists and is executable.
<code>-O file</code>	Checks if <code>file</code> exists and is owned by the current user.
<code>-G file</code>	Checks if <code>file</code> exists and the default group is the same as the current user.
<code>file1 -nt file2</code>	Checks if <code>file1</code> is newer than <code>file2</code> .
<code>file1 -ot file2</code>	Checks if <code>file1</code> is older than <code>file2</code> .

These conditions give you the ability to check files in your filesystem within your shell scripts, and they are often used in scripts that access files. Because they're used so much, let's look at each of these individually.

Checking Directories

The `-d` test checks if a specified filename exists as a directory on the system. This is usually a good thing to do if you're trying to write a file to a directory, or before you try to change to a directory location:

```
$ cat test11
#!/bin/bash
# look before you leap
if [ -d $HOME ]
then
    echo "Your HOME directory exists"
    cd $HOME
    ls -a
else
    echo "There is a problem with your HOME directory"
fi
$
$ ./test11
"Your HOME directory exists"
.      Documents      .gvfs      .pulse-cookie
```

```

Downloads      .ICEauthority      .recently-used.xbel
.aptitude      .esd_auth          .local
.sudo_as_admin_successful
.bash_history  examples          .desktop          .mozilla          Templates
.bash_logout  .fontconfig       Music             test11
.bashrc       .gconf            .nautilus         Videos
.cache       .gconfd          .openoffice.org
.xsession-errors
.config      .gksu.lock       Pictures
.xsession-errors.old
.dbus       .gnome2          .profile
Desktop    .gnome2_private Public
.dmrc     .gtk-bookmarks  .pulse
$

```

The sample code uses the `-d` test condition to see if the `$HOME` directory exists for the user. If it does, it proceeds to use the `cd` command to change to the `$HOME` directory and performs a directory listing.

Checking If an Object Exists

The `-e` comparison allows you to check if a file or directory object exists before you attempt to use it in your script:

```

$ cat test12
#!/bin/bash
# checking if a directory exists
if [ -e $HOME ]
then
    echo "OK on the directory, now to check the file"
    # checking if a file exists
    if [ -e $HOME/testing ]
    then
        # the file exists, append data to it
        echo "Appending date to existing file"
        date >> $HOME/testing
    else
        # the file does not exist, create a new file
        echo "Creating new file"
        date > $HOME/testing
    fi
else
    echo "Sorry, you do not have a HOME directory"
fi

```

```
$ ./test12
OK on the directory, now to check the file
Creating new file
$ ./test12
OK on the directory, now to check the file
Appending date to existing file
$
```

The first check uses the `-e` comparison to determine if the user has a `$HOME` directory. If so, the next `-e` comparison checks to determine if the testing file exists in the `$HOME` directory. If the file doesn't exist, the shell script uses the single greater-than redirect symbol, creating a new file with the output from the `date` command. The second time you run the shell script, it uses the double greater-than symbol, so it just appends the date output to the existing file.

Checking for a File

The `-e` comparison works for both files and directories. To be sure that the object specified is a file, you must use the `-f` comparison:

```
$ cat test13
#!/bin/bash
# check if a file
if [ -e $HOME ]
then
    echo "The object exists, is it a file?"
    if [ -f $HOME ]
    then
        echo "Yes, it is a file!"
    else
        echo "No, it is not a file!"
        if [ -f $HOME/.bash_history ]
        then
            echo "But this is a file!"
        fi
    fi
else
```

```

echo "Sorry, the object does not exist"
fi
$
$ ./test13
The object exists, is it a file?
No, it is not a file!
But this is a file!
$

```

This little script does a whole lot of checking! First, it uses the `-e` comparison to test if `$HOME` exists. If it does, it uses `-f` to test if it's a file. If it isn't a file (which of course it isn't), we use the `-f` comparison to test if it's a file, which it is.

Can You Read It?

Before trying to read data from a file, it's usually a good idea to test if you can read from the file first. You do this with the `-r` comparison:

```

$ cat test14
#!/bin/bash
# testing if you can read a file
pwfile=/etc/shadow

# first, test if the file exists, and is a file
if [ -f $pwfile ]
then
    # now test if you can read it
    if [ -r $pwfile ]
    then
        tail $pwfile
    else
        echo "Sorry, I am unable to read the $pwfile file"
    fi
else
    echo "Sorry, the file $file does not exist"
fi
$
$ ./test14
Sorry, I am unable to read the /etc/shadow file
$

```

The `/etc/shadow` file contains the encrypted passwords for system users, so it's not readable by normal users on the system. The `-r` comparison determined that I didn't have read access to the file, so the test command failed and the bash shell executed the else section of the `if-then` statement.

Checking for Empty Files

You should use `-s` comparison to check if a file is empty, especially if you're trying to remove a file. Be careful because when the `-s` comparison succeeds, it indicates that a file has data in it:

```
$ cat test15
#!/bin/bash
# testing if a file is empty
file=tl5test
touch $file

if [ -s $file ]
then
    echo "The $file file exists and has data in it"
else
    echo "The $file exists and is empty"
fi
date > $file
if [ -s $file ]
then
    echo "The $file file has data in it"
else
    echo "The $file is still empty"
fi
$
$ ./test15
The tl5test exists and is empty
The tl5test file has data in it
$
```

The `touch` command creates the file but doesn't put any data in it. After we use the `date` command and redirect the

output to the file, the `-s` comparison indicates that there's data in the file.

Checking If You Can Write to a File

The `-w` comparison determines if you have permission to write to a file:

```
$ cat test16
#!/bin/bash
# checking if a file is writeable

logfile=$HOME/tl6test
touch $logfile
chmod u-w $logfile
now='date +%Y%m%d-%H%M'

if [ -w $logfile ]
then
    echo "The program ran at: $now" > $logfile
    echo "The first attempt succeeded"
else
    echo "The first attempt failed"
fi

chmod u+w $logfile
if [ -w $logfile ]
then
    echo "The program ran at: $now" > $logfile
    echo "The second attempt succeeded"
else
    echo "The second attempt failed"
fi
$
$ ./test16
The first attempt failed
The second attempt succeeded
$ cat $HOME/tl6test
The program ran at: 20110124-1602
$
```

This is a pretty busy shell script! First, it defines a log file in your `$HOME` directory, stores the filename of it in the variable `logfile`, creates the file, and then removes the write

permission for the user, using the `chmod` command. Next, it creates the variable `now` and stores a timestamp, using the `date` command. After all of that, it checks if you have write permission to the new log file (which you just took away). Because you don't have write permission, you should see the failed message appear.

After that, the script uses the `chmod` command again to grant the write permission back for the user, and tries to write to the file again. This time the write is successful.

Checking If You Can Run a File

The `-x` comparison is a handy way to determine if you have execute permission for a specific file. While this may not be needed for most commands, if you run a lot of scripts from your shell scripts, it could come in handy:

```
$ cat test17
#!/bin/bash
# testing file execution
if [ -x test16 ]
then
    echo "You can run the script: "
    ./test16
else
    echo "Sorry, you are unable to execute the script"
fi
$
$ ./test17
You can run the script:
The first attempt failed
The second attempt succeeded
$
$ chmod u-x test16
$
$ ./test17
Sorry, you are unable to execute the script
$
```


This example shell script uses the `-x` comparison to test if you have permission to execute the `test16` script. If so, it runs the script (notice that even in a shell script, you must have the proper path to execute a script that's not located in your `$PATH`). After successfully running the `test16` script the first time, change the permissions on it and try again. This time, the `-x` comparison fails, as you don't have execute permission for the `test16` script.

Checking Ownership

The `-o` comparison allows you to easily test if you're the owner of a file:

```
$ cat test18
#!/bin/bash
# check file ownership

if [ -o /etc/passwd ]
then
    echo "You are the owner of the /etc/passwd file"
else
    echo "Sorry, you are not the owner of the /etc/passwd file"
fi
$
$ ./test18
Sorry, you are not the owner of the /etc/passwd file
$
$ su
Password:
$
# ./test18
You are the owner of the /etc/passwd file
#
```

The script uses the `-o` comparison to test if the user running the script is the owner of the `/etc/passwd` file. The first time the script is run under a normal user account, so the test fails. The second time, we used the `su` command to

become the root user, and the test succeeded.

Checking Default Group Membership

The `-G` comparison checks the default group of a file, and it succeeds if it matches the group of the default group for the user. This can be somewhat confusing because the `-G` comparison checks the default groups only and not all the groups the user belongs to. Here's an example of this:

```
$ cat test19
#!/bin/bash# check file group test

if [ -G $HOME/testing ]
then
    echo "You are in the same group as the file"
else
    echo "The file is not owned by your group"
fi
$
$ ls -l $HOME/testing
-rw-rw-r-- 1 rich rich 58 2011-01-30 15:51 /home/rich/testing
$
$ ./test19
You are in the same group as the file
$
$ chgrp sharing $HOME/testing
$
$ ./test19
The file is not owned by your group
$
```

The first time the script is run, the `$HOME/testing` file is in the `rich` group, and the `-G` comparison succeeds. Next, the group is changed to the `sharing` group, which the user is also a member of. However, the `-G` comparison failed, since it only compares the default groups, not any additional group memberships.

Checking File Date

The last set of comparisons deal with comparing the creation times of two files. This comes in handy when writing scripts to install software. Sometimes you don't want to install a file that is older than a file already installed on the system.

The `-nt` comparison determines if a file is newer than another file. If a file is newer, it will have a more recent file creation time. The `-ot` comparison determines if a file is older than another file. If the file is older, it will have an older file creation time:

```
$ cat test20
#!/bin/bash
# testing file dates

if [ ./test19 -nt ./test18 ]
then
    echo "The test19 file is newer than test18"
else
    echo "The test18 file is newer than test19"
fi
if [ ./test17 -ot ./test19 ]
then
    echo "The test17 file is older than the test19 file"
fi
$
$ ./test20
The test19 file is newer than test18
The test17 file is older than the test19 file
$
$ ls -l test17 test18 test19
-rwxrw-r-- 1 rich rich 167 2011-01-30 16:31 test17
-rwxrw-r-- 1 rich rich 185 2011-01-30 17:46 test18
-rwxrw-r-- 1 rich rich 167 2011-01-30 17:50 test19
$
```

The filepaths used in the comparisons are relative to the directory from where you run the script. This can cause

problems if the files you're checking can be moved around. Another problem is that neither of these comparisons check if the file exists first. Try this test:

```
$ cat test21
#!/bin/bash
# testing file dates

if [ ./badfile1 -nt ./badfile2 ]
then
    echo "The badfile1 file is newer than badfile2"
else
    echo "The badfile2 file is newer than badfile1"
fi
$
$ ./test21
The badfile2 file is newer than badfile1
$
```

This little example demonstrates that if the files don't exist, the `-nt` comparison just returns a failed condition. It's imperative that you ensure the files exist before trying to use them in the `-nt` or `-ot` comparison.

Compound Condition Testing

The `if-then` statement allows you to use Boolean logic to combine tests. There are two Boolean operators you can use:

- `[condition1] && [condition2]`
- `[condition1] || [condition2]`

The first Boolean operation uses the `AND` Boolean operator to combine two conditions. Both conditions must be met for the `then` section to execute.

The second Boolean operation uses the `OR` Boolean

operator to combine two conditions. If either condition evaluates to a true condition, the `then` section is executed.

```
$ cat test22
#!/bin/bash
# testing compound comparisons

if [ -d $HOME ] && [ -w $HOME/testing ]
then
    echo "The file exists and you can write to it"
else
    echo "I cannot write to the file"
fi
$
$ ./test22
I cannot write to the file
$ touch $HOME/testing
$
$ ./test22
The file exists and you can write to it
$
```

Using the `AND` Boolean operator, both of the comparisons must be met. The first comparison checks to see if the `$HOME` directory exists for the user. The second comparison checks to see if there's a file called `testing` in the user's `$HOME` directory, and if the user has write permissions for the file. If either of these comparisons fails, the `if` statement fails and the shell executes the `else` section. If both of the comparisons succeed, the `if` statement succeeds and the shell executes the `then` section.

Advanced if-then Features

There are two relatively recent additions to the bash shell that provide advanced features that you can use in `if-then` statements:

- Double parentheses for mathematical expressions
- Double square brackets for advanced string handling functions

The following sections describe each of these features in more detail.

Using Double Parentheses

The *double parentheses* command allows you to incorporate advanced mathematical formulas in your comparisons. The `test` command only allows for simple arithmetic operations in the comparison. The *double parentheses* command provides more mathematical symbols, which programmers, who have used other programming languages, may be used to using. The format of the *double parentheses* command is:

```
(( expression ))
```

The *expression* term can be any mathematical assignment or comparison expression. Besides the standard mathematical operators that the `test` command uses, [Table 11.4](#) shows the list of additional operators available for use in the *double parentheses* command.

Table 11.4 The Double Parentheses Command Symbols

Symbol	Description
<code>val++</code>	Post-increment
<code>val--</code>	Post-decrement
<code>++val</code>	Pre-increment

--val	Pre-decrement
!	Logical negation
~	Bitwise negation
**	Exponentiation
<<	Left bitwise shift
>>	Right bitwise shift
&	Bitwise Boolean AND
	Bitwise Boolean OR
&&	Logical AND
	Logical OR

You can use the double parentheses command in an `if` statement, as well as a normal command in the script for assigning values:

```
$ cat test23
#!/bin/bash
# using double parenthesis

vall=10

if (( $vall ** 2 > 90 ))
then
    (( val2 = $vall ** 2 ))
    echo "The square of $vall is $val2"
fi
$
$ ./test23
The square of 10 is 100
$
```

Notice that you don't need to escape the greater-than symbol in the expression within the double parentheses. This is yet another advanced feature provided by the double parentheses command.

Using Double Brackets

The *double bracket* command provides advanced features for string comparisons. The double bracket command format is:

```
[[ expression ]]
```

The double bracketed *expression* uses the standard string comparison used in the `test` command. However, it provides an additional feature that the `test` command doesn't, *pattern matching*.

In pattern matching, you can define a regular expression (discussed in detail in Chapter 19) that's matched against the string value:

```
$ cat test24
#!/bin/bash
# using pattern matching

if [[ $USER == r* ]]
then
    echo "Hello $USER"
else
    echo "Sorry, I do not know you"
fi
$
$ ./test24
Hello rich
$
```

The double bracket command matches the `$USER` environment variable to see if it starts with the letter `r`. If so, the comparison succeeds, and the shell executes the `then` section commands.

The case Command

Often you'll find yourself trying to evaluate the value of a variable, looking for a specific value within a set of possible values. In this scenario, you end up having to write a lengthy `if-then-else` statement, like this:

```
$ cat test25
#!/bin/bash
# looking for a possible value

if [ $USER = "rich" ]
then
    echo "Welcome $USER"
    echo "Please enjoy your visit"
elif [ $USER = barbara ]
then
    echo "Welcome $USER"
    echo "Please enjoy your visit"
elif [ $USER = testing ]
then
    echo "Special testing account"
elif [ $USER = jessica ]
then
    echo "Do not forget to logout when you're done"
else
    echo "Sorry, you are not allowed here"
fi
$
$ ./test25
Welcome rich
Please enjoy your visit
$
```

The `elif` statements continue the `if-then` checking, looking for a specific value for the single comparison variable.

Instead of having to write all of the `elif` statements to continue checking the same variable value, you can use the `case` command. The `case` command checks multiple values of a single variable in a list-oriented format:

```
case variable in
pattern1 | pattern2) commands1::
```

```
pattern3) commands2;;
*) default commands;;
esac
```

The `case` command compares the variable specified against the different patterns. If the variable matches the pattern, the shell executes the commands specified for the pattern. You can list more than one pattern on a line, using the bar operator to separate each pattern. The asterisk symbol is the catch-all for values that don't match any of the listed patterns. Here's an example of converting the `if-then-else` program to using the `case` command:

```
$ cat test26
#!/bin/bash
# using the case command

case $USER in
rich | barbara)
    echo "Welcome, $USER"
    echo "Please enjoy your visit";;
testing)
    echo "Special testing account";;
jessica)
    echo "Do not forget to log off when you're done";;
*)
    echo "Sorry, you are not allowed here";;
esac
$
$ ./test26
Welcome, rich
Please enjoy your visit
$
```

The `case` command provides a much cleaner way of specifying the various options for each possible variable value.

Summary

Structured commands allow you to alter the normal flow of execution on the shell script. The most basic structured command is the `if-then` statement. This statement allows you to evaluate a command, and perform other commands based on the outcome of the command you evaluated.

You can expand the `if-then` statement to include a set of commands the bash shell executes if the specified command fails as well. The `if-then-else` statement allows you to execute commands only if the command being evaluated returns a non-zero exit status code.

You can also link `if-then-else` statements together, using the `elif` statement. The `elif` is equivalent to using an `else if` statement, providing for additional checking if the original command that was evaluated failed.

In most scripts, instead of evaluating a command, you'll want to evaluate a condition, such as a numeric value, the contents of a string, or the status of a file or directory. The `test` command provides an easy way for you to evaluate all of these conditions. If the condition evaluates to a true condition, the `test` command produces a zero exit status code for the `if-then` statement. If the condition evaluates to a false condition, the `test` command produces a non-zero exit status code for the `if-then` statement.

The square bracket is a special bash command that is a synonym for the `test` command. You can enclose a test condition in square brackets in the `if-then` statement to test for numeric, string, and file conditions.

The double parentheses command allows you to perform

advanced mathematical evaluations using additional operators, and the double square bracket command allows you to perform advanced string pattern-matching evaluations.

Finally, the chapter discussed the `case` command, which is a shorthand way of performing multiple `if-then-else` commands, checking the value of a single variable against a list of values.

The next chapter continues the discussion of structured commands by examining the shell looping commands. The `for` and `while` commands allow you to create loops that iterate through commands for a given period of time.

Chapter 12

More Structured Commands

In This Chapter

- Looping with the `for` statement
- Iterating with the `until` statement
- Using the `while` statement
- Combining loops
- Redirecting loop output

In the previous chapter, you saw how to manipulate the flow of a shell script program by checking the output of commands and the values of variables. In this chapter, we'll continue to look at structured commands that control the flow of your shell scripts. You'll see how you can perform repeating processes, commands that can loop through a set of commands until an indicated condition has been met. This chapter discusses and demonstrates the `for`, `while`, and `until` bash shell looping commands.

The `for` Command

Iterating through a series of commands is a common programming practice. Often you need to repeat a set of commands until a specific condition has been met, such as processing all of the files in a directory, all of the users on a system, or all of the lines in a text file.

The bash shell provides the `for` command to allow you to create a loop that iterates through a series of values. Each iteration performs a defined set of commands using one of the values in the series. The following is the basic format of the bash shell `for` command:

```
for var in list
do
  commands
done
```

You supply the series of values used in the iterations in the *list* parameter. There are several different ways that you can specify the values in the list.

In each iteration, the variable *var* contains the current value in the list. The first iteration uses the first item in the list, the second iteration the second item, and so on until all of the items in the list have been used.

The *commands* entered between the `do` and `done` statements can be one or more standard bash shell commands. Within the commands, the `$var` variable contains the current list item value for the iteration.

Note

If you prefer, you can include the `do` statement on the same line as the `for` statement, but you must

separate it from the list items using a semicolon: `for`

```
var in list; do.
```

We mentioned that there are several different ways to specify the values in the list. The following sections show the various ways to do that.

Reading Values in a List

The most basic use of the `for` command is to iterate through a list of values defined within the `for` command itself:

```
$ cat test1
#!/bin/bash
# basic for command

for test in Alabama Alaska Arizona Arkansas California Colorado
do
    echo The next state is $test
done
$ ./test1
The next state is Alabama
The next state is Alaska
The next state is Arizona
The next state is Arkansas
The next state is California
The next state is Colorado
$
```

Each time the `for` command iterates through the list of values provided, it assigns the `$test` variable the next value in the list. The `$test` variable can be used just like any other script variable within the `for` command statements. After the last iteration, the `$test` variable remains valid throughout the remainder of the shell script. It retains the last iteration value (unless you change its value):

```
$ cat test1b
#!/bin/bash
```

```
# testing the for variable after the looping

for test in Alabama Alaska Arizona Arkansas California Colorado
do
    echo "The next state is $test"
done
echo "The last state we visited was $test"
test=Connecticut
echo "Wait, now we're visiting $test"
$ ./test1b
The next state is Alabama
The next state is Alaska
The next state is Arizona
The next state is Arkansas
The next state is California
The next state is Colorado
The last state we visited was Colorado
Wait, now we're visiting Connecticut
$
```

The `$test` variable retained its value and also allowed us to change the value and use it outside of the `for` command loop, as any other variable would.

Reading Complex Values in a List

Things aren't always as easy as they seem with the `for` loop. There are times when you run into data that causes problems. Here's a classic example of what can cause problems for shell script programmers:

```
$ cat badtest1
#!/bin/bash
# another example of how not to use the for command

for test in I don't know if this'll work
do
    echo "word:$test"
done
$ ./badtest1
word:I
word:dont know if thisll
word:work
$
```


Ouch, that hurts. The shell saw the single quotation marks within the list values and attempted to use them to define a single data value, and it really messed things up in the process.

There are two ways to solve this problem:

- Use the escape character (the backslash) to escape the single quotation mark.
- Use double quotation marks to define the values that use single quotation marks.

Neither solution is all that fantastic, but each one does help solve the problem:

```
$ cat test2
#!/bin/bash
# another example of how not to use the for command

for test in I don't know if "this'll" work
do
    echo "word:$test"
done
$ ./test2
word:I
word:don't
word:know
word:if
word:this'll
word:work
$
```

In the first problem value, you added the backslash character to escape the single quotation mark in the `don't` value. In the second problem value, you enclosed the `this'll` value in double quotation marks. Both methods worked fine to distinguish the value.

Yet another problem you may run into is multi-word values. Remember that the `for` loop assumes that each

value is separated with a space. If you have data values that contain spaces, you'll run into yet another problem:

```
$ cat badtest2
#!/bin/bash
# another example of how not to use the for command

for test in Nevada New Hampshire New Mexico New York North Carolina
do
    echo "Now going to $test"
done
$ ./badtest1
Now going to Nevada
Now going to New
Now going to Hampshire
Now going to New
Now going to Mexico
Now going to New
Now going to York
Now going to North
Now going to Carolina
$
```

Oops, that's not exactly what we wanted. The `for` command separates each value in the list with a space. If there are spaces in the individual data values, you must accommodate them using double quotation marks:

```
$ cat test3
#!/bin/bash
# an example of how to properly define values

for test in Nevada "New Hampshire" "New Mexico" "New York"
do
    echo "Now going to $test"
done
$ ./test3
Now going to Nevada
Now going to New Hampshire
Now going to New Mexico
Now going to New York
$
```

Now the `for` command can properly distinguish between the different values. Also, notice that when you use double

quotation marks around a value, the shell doesn't include the quotation marks as part of the value.

Reading a List from a Variable

Often what happens in a shell script is that you accumulate a list of values stored in a variable and then need to iterate through the list. You can do this using the `for` command as well:

```
$ cat test4
#!/bin/bash
# using a variable to hold the list

list="Alabama Alaska Arizona Arkansas Colorado"
list=$list" Connecticut"

for state in $list
do
    echo "Have you ever visited $state?"
done
$ ./test4
Have you ever visited Alabama?
Have you ever visited Alaska?
Have you ever visited Arizona?
Have you ever visited Arkansas?
Have you ever visited Colorado?
Have you ever visited Connecticut?
$
```

The `$list` variable contains the standard text list of values to use for the iterations. Notice that the code also uses another assignment statement to add (or concatenate) an item to the existing list contained in the `$list` variable. This is a common method for adding text to the end of an existing text string stored in a variable.

Reading Values from a Command

Yet another way to generate values for use in the list is to use the output of a command. You use the backtick characters to execute any command that produces output, and then use the output of the command in the `for` command:

```
$ cat test5
#!/bin/bash
# reading values from a file

file="states"

for state in `cat $file`
do
    echo "Visit beautiful $state"
done
$ cat states
Alabama
Alaska
Arizona
Arkansas
Colorado
Connecticut
Delaware
Florida
Georgia
$ ./test5
Visit beautiful Alabama
Visit beautiful Alaska
Visit beautiful Arizona
Visit beautiful Arkansas
Visit beautiful Colorado
Visit beautiful Connecticut
Visit beautiful Delaware
Visit beautiful Florida
Visit beautiful Georgia
$
```

This example uses the `cat` command in tick marks to display the contents of the file `states`. You'll notice that the `states` file includes each state on a separate line, not separated by spaces. The `for` command still iterates through the output of the `cat` command one line at a time,

assuming that each state is on a separate line. However, this doesn't solve the problem of having spaces in data. If you list a state with a space in it, the `for` command will still take each word as a separate value. There's a reason for this, which we look at in the next section.

Note

The `test5` code example assigned the filename to the variable using just the filename without a path. This requires that the file be in the same directory as the script. If this isn't the case, you'll need to use a full pathname (either absolute or relative) to reference the file location.

Changing the Field Separator

The cause of this problem is the special environment variable `IFS`, called the *internal field separator*. The `IFS` environment variable defines a list of characters the bash shell uses as field separators. By default, the bash shell considers the following characters as field separators:

- A space
- A tab
- A newline

If the bash shell sees any of these characters in the data, it will assume that you're starting a new data field in the list. When working with data that can contain spaces (such as filenames), this can be annoying, as you saw in the

previous script example.

To solve this problem, you can temporarily change the IFS environment variable values in your shell script to restrict the characters the bash shell recognizes as field separators. However, there is somewhat of an odd way of doing this. For example, if you want to change the IFS value to recognize only the newline character, you need to do this:

```
IFS=$'\n'
```

Adding this statement to your script tells the bash shell to ignore spaces and tabs in data values. Applying this to the previous script yields the following:

```
$ cat test5b
#!/bin/bash
# reading values from a file

file="states"

IFS=$'\n'
for state in `cat $file`
do
    echo "Visit beautiful $state"
done
$ ./test5b
Visit beautiful Alabama
Visit beautiful Alaska
Visit beautiful Arizona
Visit beautiful Arkansas
Visit beautiful Colorado
Visit beautiful Connecticut
Visit beautiful Delaware
Visit beautiful Florida
Visit beautiful Georgia
Visit beautiful New York
Visit beautiful New Hampshire
Visit beautiful North Carolina
$
```

Now the shell script is able to use values in the list that contain spaces.

Caution

When working on long scripts, it's possible to change the IFS value in one place, and then forget about it and assume the default value elsewhere in the script. A safe practice to get into is to save the original IFS value before changing it and then restore it when you're done.

This technique can be coded like this:

```
IFS.OLD=$IFS
IFS=$'\n'
<use the new IFS value in code>
IFS=$IFS.OLD
```

This ensures that the IFS value is returned to the default value for future operations within the script.

There are other excellent applications of the IFS environment variable. Say that you want to iterate through values in a file that are separated by a colon (such as in the `/etc/passwd` file). All you need to do is set the IFS value to a colon:

```
IFS=:
```

If you want to specify more than one IFS character, just string them together on the assignment line:

```
IFS=$'\n';;"
```

This assignment uses the newline, colon, semicolon, and double quotation mark characters as field separators. There's no limit to how you can parse your data using the IFS characters.

Reading a Directory Using Wildcards

Finally, you can use the `for` command to automatically iterate through a directory of files. To do this, you must use a wildcard character in the file or pathname. This forces the shell to use *file globbing*. File globbing is the process of producing file or path names that match a specified wildcard character.

This feature is great for processing files in a directory when you don't know all of the filenames:

```
$ cat test6
#!/bin/bash
# iterate through all the files in a directory

for file in /home/rich/test/*
do

    if [ -d "$file" ]
    then
        echo "$file is a directory"
    elif [ -f "$file" ]
    then
        echo "$file is a file"
    fi
done

$ ./test6
/home/rich/test/dirl is a directory
/home/rich/test/myprog.c is a file
/home/rich/test/myprog is a file
/home/rich/test/myscript is a file
/home/rich/test/newdir is a directory
/home/rich/test/newfile is a file
/home/rich/test/newfile2 is a file
/home/rich/test/testdir is a directory
/home/rich/test/testing is a file
/home/rich/test/testprog is a file
/home/rich/test/testprog.c is a file
$
```

The `for` command iterates through the results of the `/home/rich/test/*` listing. The code tests each entry using the `test` command (using the square bracket method) to see if

it's a directory, using the `-d` parameter, or a file, using the `-f` parameter (See Chapter 11).

Notice in this example that we did something different in the `if` statement tests:

```
if [ -d "$file" ]
```

In Linux, it's perfectly legal to have directory and filenames that contain spaces. To accommodate that, you should enclose the `$file` variable in double quotation marks. If you don't, you'll get an error if you run into a directory or filename that contains spaces:

```
./test6: line 6: [: too many arguments
./test6: line 9: [: too many arguments
```

The `bash` shell interprets the additional words as arguments within the `test` command, causing an error.

You can also combine both the directory search method and the list method in the same `for` statement by listing a series of directory wildcards in the `for` command:

```
$ cat test7
#!/bin/bash
# iterating through multiple directories

for file in /home/rich/.b* /home/rich/badtest
do
    if [ -d "$file" ]
    then
        echo "$file is a directory"
    elif [ -f "$file" ]
    then
        echo "$file is a file"
    else
        echo "$file doesn't exist"
    fi
done
$ ./test7
/home/rich/.backup.timestamp is a file
/home/rich/.bash_history is a file
/home/rich/.bash_logout is a file
```

```
/home/rich/.bash_profile is a file
/home/rich/.bashrc is a file
/home/rich/badtest doesn't exist
$
```

The `for` statement first uses file globbing to iterate through the list of files that result from the wildcard character; then it iterates through the next file in the list. You can combine any number of wildcard entries in the list to iterate through.

Caution

Notice that you can enter anything in the list data. Even if the file or directory doesn't exist, the `for` statement attempts to process whatever you place in the list. This can be a problem when working with files and directories. You have no way of knowing if you're trying to iterate through a nonexistent directory. It's always a good idea to test each file or directory before trying to process it.

The C-Style for Command

If you've done any programming using the C programming language, you're probably surprised by the way the bash shell uses the `for` command. In the C language, a `for` loop normally defines a variable, which it then alters automatically during each iteration. Typically, programmers use this variable as a counter and either increment or decrement the counter by one in each iteration. The bash `for` command can also provide this functionality. This

section shows you how you can use a C-style `for` command in a bash shell script.

The C Language for Command

The C language `for` command has a specific method for specifying a variable, a condition that must remain true for the iterations to continue, and a method for altering the variable for each iteration. When the specified condition becomes false, the `for` loop stops. The condition equation is defined using standard mathematical symbols. For example, consider the following C language code:

```
for (i = 0; i < 10; i++)
{
    printf("The next number is %d\n", i);
}
```

This code produces a simple iteration loop, where the variable `i` is used as a counter. The first section assigns a default value to the variable. The middle section defines the condition under which the loop will iterate. When the defined condition becomes false, the `for` loop stops iterations. The last section defines the iteration process. After each iteration, the expression defined in the last section is executed. In this example, the `i` variable is incremented by one after each iteration.

The bash shell also supports a version of the `for` loop that looks similar to the C-style `for` loop, although it does have some subtle differences, including a couple of things that will confuse shell script programmers. Here's the basic format of the C-style bash `for` loop:

```
for ((  
variable  
assignment ;  
condition ;  
iteration  
process ))
```

The format of the C-style `for` loop can be confusing for bash shell script programmers, as it uses C-style variable references instead of the shell-style variable references. Here's what a C-style `for` command looks like:

```
for (( a = 1; a < 10; a++ ))
```

Notice that there are a couple of things that don't follow the standard bash shell `for` method:

- The assignment of the variable value can contain spaces.
- The variable in the condition isn't preceded with a dollar sign.
- The equation for the iteration process doesn't use the `expr` command format.

The shell developers created this format to more closely resemble the C-style `for` command. While this is great for C programmers, it can throw even expert shell programmers into a tizzy. Be careful when using the C-style `for` loop in your scripts.

Here's an example of using the C-style `for` command in a bash shell program:

```
$ cat test8  
#!/bin/bash  
# testing the C-style for loop  
  
for (( i=1; i <= 10; i++ ))  
do  
    echo "The next number is $i"
```

```
done
$ ./test8
The next number is 1
The next number is 2
The next number is 3
The next number is 4
The next number is 5
The next number is 6
The next number is 7
The next number is 8
The next number is 9
The next number is 10
$
```

The `for` loop iterates through the commands using the variable defined in the `for` loop (the letter *i* in this example). In each iteration, the `$i` variable contains the value assigned in the `for` loop. After each iteration, the loop iteration process is applied to the variable, which in this example, increments the variable by one.

Using Multiple Variables

The C-style `for` command also allows you to use multiple variables for the iteration. The loop handles each variable separately, allowing you to define a different iteration process for each variable. While you can have multiple variables, you can only define one condition in the `for` loop:

```
$ cat test9
#!/bin/bash
# multiple variables

for (( a=1, b=10; a <= 10; a++, b-- ))
do
    echo "$a - $b"
done
$ ./test9
1 - 10
2 - 9
3 - 8
```

```
4 - 7
5 - 6
6 - 5
7 - 4
8 - 3
9 - 2
10 - 1
$
```

The `a` and `b` variables are each initialized with different values and different iteration processes are defined. While the loop increases the `a` variable, it decreases the `b` variable for each iteration.

The while Command

The `while` command is somewhat of a cross between the `if-then` statement and the `for` loop. The `while` command allows you to define a command to test and then loop through a set of commands for as long as the defined test command returns a zero exit status. It tests the `test` command at the start of each iteration. When the `test` command returns a non-zero exit status, the `while` command stops executing the set of commands.

Basic while Format

The format of the `while` command is:

```
while
test
command
do

other
commands
```

done

The `test command` defined in the `while` command is the exact same format as in `if-then` statements (see Chapter 11). As in the `if-then` statement, you can use any normal bash shell command, or you can use the `test` command to test for conditions, such as variable values.

The key to the `while` command is that the exit status of the `test command` specified must change, based on the commands run during the loop. If the exit status never changes, the `while` loop will get stuck in an infinite loop.

The most common use of the `test command` is to use brackets to check a value of a shell variable that's used in the loop commands:

```
$ cat test10
#!/bin/bash
# while command test

var1=10
while [ $var1 -gt 0 ]
do
    echo $var1
    var1=$(( $var1 - 1 ))
done
$ ./test10
10
9
8
7
6
5
4
3
2
1
$
```

The `while` command defines the test condition to check for each iteration:

```
while [ $var1 -gt 0 ]
```

As long as the test condition is true, the `while` command continues to loop through the commands defined. Within the commands, the variable used in the test condition must be modified, or else you'll have an infinite loop. In this example, we use shell arithmetic to decrease the variable value by one:

```
var1=$(( var1 - 1 )
```

The `while` loop stops when the test condition is no longer true.

Using Multiple Test Commands

In somewhat of an odd situation, the `while` command allows you to define multiple test commands on the `while` statement line. Only the exit status of the last test command is used to determine when the loop stops. This can cause some interesting results if you're not careful. Here's an example of what we mean:

```
$ cat test11
#!/bin/bash
# testing a multicommand while loop
```

```
var1=10
```

```
while echo $var1
      [ $var1 -ge 0 ]
do
    echo "This is inside the loop"
    var1=$(( var1 - 1 )
```

```
done
```

```
$ ./test11
```

```
10
```

```
This is inside the loop
```

```
9
```



```
This is inside the loop
8
This is inside the loop
7
This is inside the loop
6
This is inside the loop
5
This is inside the loop
4
This is inside the loop
3
This is inside the loop
2
This is inside the loop
1
This is inside the loop
0
This is inside the loop
-1
$
```

Pay close attention to what happened in this example. There were two test commands defined in the `while` statement:

```
while echo $var1
[ $var1 -ge 0 ]
```

The first test simply displays the current value of the `var1` variable. The second test uses brackets to determine the value of the `var1` variable. Inside the loop, an `echo` statement displays a simple message, indicating that the loop was processed. Notice when you run the example how the output ends:

```
This is inside the loop
-1
$
```

The `while` loop executed the `echo` statement when the `var1` variable was equal to zero, and then decreased the `var1`

variable value. Next, the test commands were executed for the next iteration. The `echo` test command was executed, displaying the value of the `var1` variable, which is now less than zero. It's not until the shell executes the `test` test command that the `while` loop terminates.

This demonstrates that in a multi-command `while` statement, all of the test commands are executed in each iteration, including the last iteration when the last test command fails. Be careful of this. Another thing to be careful of is how you specify the multiple test commands. Note that each test command is on a separate line!

The until Command

The `until` command works in exactly the opposite way from the `while` command. The `until` command requires that you specify a test command that normally produces a non-zero exit status. As long as the exit status of the test command is non-zero, the bash shell executes the commands listed in the loop. Once the test command returns a zero exit status, the loop stops.

As you would expect, the format of the `until` command is:

```
until test
commands
do
    other commands
done
```

Similar to the `while` command, you can have more than one *test command* in the `until` command statement. Only

the exit status of the last command determines if the bash shell executes the *other commands* defined.

The following is an example of using the `until` command:

```
$ cat test12
#!/bin/bash
# using the until command

var1=100

until [ $var1 -eq 0 ]
do
    echo $var1
    var1=$(( $var1 - 25 ])
done
$ ./test12
100
75
50
25
$
```

This example tests the `var1` variable to determine when the `until` loop should stop. As soon as the value of the variable is equal to zero, the `until` command stops the loop. The same caution as for the `while` command applies when you use multiple test commands with the `until` command:

```
$ cat test13
#!/bin/bash
# using the until command

var1=100

until echo $var1
    [ $var1 -eq 0 ]
do
    echo Inside the loop: $var1
    var1=$(( $var1 - 25 ])
done
$ ./test13
100
Inside the loop: 100
75
Inside the loop: 75
```

```
50
Inside the loop: 50
25
Inside the loop: 25
0
$
```

The shell executes the test commands specified and stops only when the last command is true.

Nesting Loops

A loop statement can use any other type of command within the loop, including other loop commands. This is called a *nested loop*. Care should be taken when using nested loops, as you're performing an iteration within an iteration, which multiplies the number of times commands are being run. Not paying close attention to this can cause problems in your scripts.

Here's a simple example of nesting a `for` loop inside another `for` loop:

```
$ cat test14
#!/bin/bash
# nesting for loops

for (( a = 1; a <= 3; a++ ))
do
    echo "Starting loop $a:"
    for (( b = 1; b <= 3; b++ ))
    do
        echo "    Inside loop: $b"
    done
done
$ ./test14
Starting loop 1:
    Inside loop: 1
    Inside loop: 2
    Inside loop: 3
```

```
Starting loop 2:
    Inside loop: 1
    Inside loop: 2
    Inside loop: 3
Starting loop 3:
    Inside loop: 1
    Inside loop: 2
    Inside loop: 3
$
```

The nested loop (also called the *inner loop*) iterates through its values for each iteration of the outer loop. Notice that there's no difference between the `do` and `done` commands for the two loops. The bash shell knows when the first `done` command is executed that it refers to the inner loop and not the outer loop.

The same applies when you mix loop commands, such as placing a `for` loop inside a `while` loop:

```
$ cat test15
#!/bin/bash
# placing a for loop inside a while loop

var1=5

while [ $var1 -ge 0 ]
do
    echo "Outer loop: $var1"
    for (( var2 = 1; $var2 < 3; var2++ ))
    do
        var3=$(( $var1 * $var2 ))
        echo "    Inner loop: $var1 * $var2 = $var3"
    done
    var1=$(( $var1 - 1 ))
done
$ ./test15
Outer loop: 5
    Inner loop: 5 * 1 = 5
    Inner loop: 5 * 2 = 10
Outer loop: 4
    Inner loop: 4 * 1 = 4
    Inner loop: 4 * 2 = 8
Outer loop: 3
    Inner loop: 3 * 1 = 3
    Inner loop: 3 * 2 = 6
```

```

Outer loop: 2
    Inner loop: 2 * 1 = 2
    Inner loop: 2 * 2 = 4
Outer loop: 1
    Inner loop: 1 * 1 = 1
    Inner loop: 1 * 2 = 2
Outer loop: 0
    Inner loop: 0 * 1 = 0
    Inner loop: 0 * 2 = 0
$

```

Again, the shell was able to distinguish between the `do` and `done` commands of the inner `for` loop from the same commands in the outer `while` loop.

If you really want to test your brain, you can even combine `until` and `while` loops:

```

$ cat test16
#!/bin/bash
# using until and while loops

var1=3

until [ $var1 -eq 0 ]
do
    echo "Outer loop: $var1"
    var2=1
    while [ $var2 -lt 5 ]
    do
        var3=`echo "scale=4; $var1 / $var2" | bc`
        echo "    Inner loop: $var1 / $var2 = $var3"
        var2=$(( $var2 + 1 ))
    done
    var1=$(( $var1 - 1 ))
done
$ ./test16
Outer loop: 3
    Inner loop: 3 / 1 = 3.0000
    Inner loop: 3 / 2 = 1.5000
    Inner loop: 3 / 3 = 1.0000
    Inner loop: 3 / 4 = .7500
Outer loop: 2
    Inner loop: 2 / 1 = 2.0000
    Inner loop: 2 / 2 = 1.0000
    Inner loop: 2 / 3 = .6666
    Inner loop: 2 / 4 = .5000
Outer loop: 1

```

```
Inner loop: 1 / 1 = 1.0000  
Inner loop: 1 / 2 = .5000  
Inner loop: 1 / 3 = .3333  
Inner loop: 1 / 4 = .2500
```

\$

The outer `until` loop starts with a value of 3 and continues until the value equals 0. The inner `while` loop starts with a value of 1 and continues as long as the value is less than 5. Each loop must change the value used in the test condition, or the loop will get stuck infinitely.

Looping on File Data

Often, you must iterate through items stored inside a file. This requires combining two of the techniques covered:

- Using nested loops
- Changing the `IFS` environment variable

By changing the `IFS` environment variable, you can force the `for` command to handle each line in the file as a separate item for processing, even if the data contains spaces. Once you've extracted an individual line in the file, you may have to loop again to extract data contained within it.

The classic example of this is processing data in the `/etc/passwd` file. This requires that you iterate through the `/etc/passwd` file line by line and then change the `IFS` variable value to a colon so that you can separate out the individual components in each line.

The following is an example of doing just that:

```
#!/bin/bash
# changing the IFS value

IFS.OLD=$IFS
IFS=$'\n'
for entry in `cat /etc/passwd`
do
    echo "Values in $entry -"
    IFS=:
    for value in $entry
    do
        echo "    $value"
    done
done
$
```

This script uses two different `IFS` values to parse the data. The first `IFS` value parses the individual lines in the `/etc/passwd` file. The inner `for` loop next changes the `IFS` value to the colon, which allows you to parse the individual values within the `/etc/passwd` lines.

When you run this script, you'll get output something like this:

```
Values in rich:x:501:501:Rich Blum:/home/rich:/bin/bash -
    rich
    x
    501
    501
    Rich Blum
    /home/rich
    /bin/bash
Values in katie:x:502:502:Katie Blum:/home/katie:/bin/bash -
    katie
    x
    506
    509
    Katie Blum
    /home/katie
    /bin/bash
```

The inner loop parses each individual value in the `/etc/passwd` entry. This is also a great way to process comma-separated data, a common way to import

Controlling the Loop

You might be tempted to think that once you start a loop, you're stuck until the loop finishes all of its iterations. This is not true. There are a couple of commands that help us control what happens inside of a loop:

- The `break` command
- The `continue` command

Each command has a different use in how to control the operation of a loop. The following sections describe how you can use these commands to control the operation of your loops.

The `break` Command

The `break` command is a simple way to escape out of a loop in progress. You can use the `break` command to exit out of any type of loop, including `while` and `until` loops.

There are several situations in which you can use the `break` command. This section shows each of these methods.

Breaking Out of a Single Loop

When the shell executes a `break` command, it attempts to break out of the loop that's currently processing:

```
$ cat test17  
#!/bin/bash
```

```
# breaking out of a for loop

for var1 in 1 2 3 4 5 6 7 8 9 10
do
    if [ $var1 -eq 5 ]
    then
        break
    fi
    echo "Iteration number: $var1"
done
echo "The for loop is completed"
$ ./test17
Iteration number: 1
Iteration number: 2
Iteration number: 3
Iteration number: 4
The for loop is completed
$
```

The `for` loop should normally have iterated through all of the values specified in the list. However, when the `if-then` condition was satisfied, the shell executed the `break` command, which stopped the `for` loop.

This technique also works for `while` and `until` loops:

```
$ cat test18
#!/bin/bash
# breaking out of a while loop

var1=1

while [ $var1 -lt 10 ]
do
    if [ $var1 -eq 5 ]
    then
        break
    fi
    echo "Iteration: $var1"
    var1=$(( $var1 + 1 ))
done
echo "The while loop is completed"
$ ./test18
Iteration: 1
Iteration: 2
Iteration: 3
Iteration: 4
The while loop is completed
```

\$

The `while` loop terminated when the `if-then` condition was met, executing the `break` command.

Breaking Out of an Inner Loop

When you're working with multiple loops, the `break` command automatically terminates the innermost loop you're in:

```
$ cat test19
#!/bin/bash
# breaking out of an inner loop

for (( a = 1; a < 4; a++ ))
do
    echo "Outer loop: $a"
    for (( b = 1; b < 100; b++ ))
    do
        if [ $b -eq 5 ]
        then
            break
        fi
        echo "    Inner loop: $b"
    done
done
$ ./test19
Outer loop: 1
    Inner loop: 1
    Inner loop: 2
    Inner loop: 3
    Inner loop: 4
Outer loop: 2
    Inner loop: 1
    Inner loop: 2
    Inner loop: 3
    Inner loop: 4
Outer loop: 3
    Inner loop: 1
    Inner loop: 2
    Inner loop: 3
    Inner loop: 4
$
```

The `for` statement in the inner loop specifies to iterate until the `b` variable is equal to 100. However, the `if-then` statement in the inner loop specifies that when the `b` variable value is equal to five, the `break` command is executed. Notice that even though the inner loop is terminated with the `break` command, the outer loop continues working as specified.

Breaking Out of an Outer Loop

There may be times when you're in an inner loop but need to stop the outer loop. The `break` command includes a single command line parameter value:

```
break
n
```

where *n* indicates the level of the loop to break out of. By default, *n* is 1, indicating to break out of the current loop. If you set *n* to a value of 2, the `break` command will stop the next level of the outer loop:

```
$ cat test20
#!/bin/bash
# breaking out of an outer loop

for (( a = 1; a < 4; a++ ))
do
    echo "Outer loop: $a"
    for (( b = 1; b < 100; b++ ))
    do
        if [ $b -gt 4 ]
        then
            break 2
        fi
        echo "    Inner loop: $b"
    done
done
$ ./test20
Outer loop: 1
```

```
Inner loop: 1
Inner loop: 2
Inner loop: 3
Inner loop: 4
```

```
$
```

Now when the shell executes the `break` command, the outer loop stops.

The continue Command

The `continue` command is a way to prematurely stop processing commands inside of a loop but not terminate the loop completely. This allows you to set conditions within a loop where the shell won't execute commands. Here's a simple example of using the `continue` command in a `for` loop:

```
$ cat test21
#!/bin/bash
# using the continue command

for (( var1 = 1; var1 < 15; var1++ ))
do
    if [ $var1 -gt 5 ] && [ $var1 -lt 10 ]
    then
        continue
    fi
    echo "Iteration number: $var1"
done
$ ./test21
Iteration number: 1
Iteration number: 2
Iteration number: 3
Iteration number: 4
Iteration number: 5
Iteration number: 10
Iteration number: 11
Iteration number: 12
Iteration number: 13
Iteration number: 14
$
```

When the conditions of the `if-then` statement are met (the

value is greater than five and less than 10), the shell executes the `continue` command, which skips the rest of the commands in the loop, but keeps the loop going. When the `if-then` condition is no longer met, things return back to normal.

You can use the `continue` command in `while` and `until` loops, but be extremely careful with what you're doing. Remember, when the shell executes the `continue` command, it skips the remaining commands. If you're incrementing your test condition variable in one of those conditions, bad things will happen:

```
$ cat badtest3
#!/bin/bash
# improperly using the continue command in a while loop

var1=0

while echo "while iteration: $var1"
  [ $var1 -lt 15 ]
do
  if [ $var1 -gt 5 ] && [ $var1 -lt 10 ]
  then
    continue
  fi
  echo "    Inside iteration number: $var1"
  var1=$(( $var1 + 1 ))
done
$ ./badtest3 | more
while iteration: 0
    Inside iteration number: 0
while iteration: 1
    Inside iteration number: 1
while iteration: 2
    Inside iteration number: 2
while iteration: 3
    Inside iteration number: 3
while iteration: 4
    Inside iteration number: 4
while iteration: 5
    Inside iteration number: 5
while iteration: 6
```



```

if [ $a -gt 2 ] && [ $a -lt 4 ]
then
    continue 2
fi
var3=$(( $a * $b ))
echo "    The result of $a * $b is $var3"
done
done
$ ./test22
Iteration 1:
    The result of 1 * 1 is 1
    The result of 1 * 2 is 2
Iteration 2:
    The result of 2 * 1 is 2
    The result of 2 * 2 is 4
Iteration 3:
Iteration 4:
    The result of 4 * 1 is 4
    The result of 4 * 2 is 8
Iteration 5:
    The result of 5 * 1 is 5
    The result of 5 * 2 is 10
$

```

The if-then statement:

```

if [ $a -gt 2 ] && [ $a -lt 4 ]
then
    continue 2
fi

```

uses the `continue` command to stop processing the commands inside the loop but continue the outer loop. Notice in the script output that the iteration for the value 3 doesn't process any inner loop statements, as the `continue` command stopped the processing, but continues with the outer loop processing.

Processing the Output of a Loop

Finally, you can either pipe or redirect the output of a loop within your shell script. You do this by adding the processing command to the end of the `done` command:

```
for file in /home/rich*
do
    if [ -d "$file" ]
    then
        echo "$file is a directory"
    elif
        echo "$file is a file"
    fi
done > output.txt
```

Instead of displaying the results on the monitor, the shell redirects the results of the `for` command to the file `output.txt`.

Consider the following example of redirecting the output of a `for` command to a file:

```
$ cat test23
#!/bin/bash
# redirecting the for output to a file

for (( a = 1; a < 10; a++ ))
do
    echo "The number is $a"
done > test23.txt
echo "The command is finished."
$ ./test23
The command is finished.
$ cat test23.txt
The number is 1
The number is 2
The number is 3
The number is 4
The number is 5
The number is 6
The number is 7
The number is 8
The number is 9
$
```

The shell creates the file `test23.txt` and redirects the output of the `for` command only to the file. The shell displays

the `echo` statement after the `for` command just as normal.

This same technique also works for piping the output of a loop to another command:

```
$ cat test24
#!/bin/bash
# piping a loop to another command

for state in "North Dakota" Connecticut Illinois Alabama Tennessee
do
    echo "$state is the next place to go"
done | sort
echo "This completes our travels"
$ ./test24
Alabama is the next place to go
Connecticut is the next place to go
Illinois is the next place to go
North Dakota is the next place to go
Tennessee is the next place to go
This completes our travels
$
```

The state values aren't listed in any particular order in the `for` command list. The output of the `for` command is piped to the `sort` command, which will change the order of the `for` command output. Running the script indeed shows that the output was properly sorted within the script.

Summary

Looping is an integral part of programming. The bash shell provides three different looping commands that you can use in your scripts. The `for` command allows you to iterate through a list of values, either supplied within the command line, contained in a variable, or obtained by using file globbing, to extract file and directory names from a

wildcard character.

The `while` command provides a method to loop based on the condition of a command, using either ordinary commands or the test command, which allows you to test conditions of variables. As long as the command (or condition) produces a zero exit status, the `while` loop will continue to iterate through the specified set of commands.

The `until` command also provides a method to iterate through commands, but it bases its iterations on a command (or condition) producing a non-zero exit status. This feature allows you to set a condition that must be met before the iteration stops.

You can combine loops in shell scripts, producing multiple layers of loops. The bash shell provides the `continue` and `break` commands, which allow you to alter the flow of the normal loop process based on different values within the loop.

The bash shell also allows you to use standard command redirection and piping to alter the output of a loop. You can use redirection to redirect the output of a loop to a file or piping to redirect the output of a loop to another command. This provides a wealth of features with which you can control your shell script execution.

The next chapter discusses how to interact with your shell script user. Often, shell scripts aren't completely self-contained. They require some sort of external data that must be supplied at the time you run them. The next chapter discusses different methods with which you can provide

real-time data to your shell scripts for processing.

Chapter 13

Handling User Input

In This Chapter

- Command line parameters
- Special parameter variables
- Being shifty
- Working with options
- Standardizing options
- Getting user input

So far you've seen how to write scripts that interact with data, variables, and files on the Linux system. Sometimes, you need to write a script that has to interact with the person running the script. The bash shell provides a few different methods for retrieving data from people, including command line parameters (data values added after the command), command line options (single-letter values that modify the behavior of the command), and the capability to read input directly from the keyboard. This chapter discusses how to incorporate these different methods into your bash shell scripts to obtain data from the person running your script.

Command Line Parameters

The most basic method of passing data to your shell script is to use *command line parameters*. Command line parameters allow you to add data values to the command line when you execute the script:

```
$ ./addem 10 30
```

This example passes two command line parameters (10 and 30) to the script `addem`. The script handles the command line parameters using special variables. The following sections describe how to use command line parameters in your bash shell scripts.

Reading Parameters

The bash shell assigns special variables, called *positional parameters*, to all of the parameters entered in a command line. This even includes the name of the program the shell executes. The positional parameter variables are standard numbers, with `$0` being the name of the program, `$1` being the first parameter, `$2` being the second parameter, and so on, up to `$9` for the ninth parameter.

Here's a simple example of using one command line parameter in a shell script:

```
$ cat test1
#!/bin/bash
# using one command line parameter

factorial=1
```

```

for (( number = 1; number <= $1 ; number++ ))
do
    factorial=$(( $factorial * $number )
done
echo The factorial of $1 is $factorial
$
$ ./test1 5
The factorial of 5 is 120
$

```

You can use the `$1` variable just like any other variable in the shell script. The shell script automatically assigns the value from the command line parameter to the variable; you don't need to do anything with it.

If you need to enter more command line parameters, each parameter must be separated by a space on the command line:

```

$ cat test2
#!/bin/bash
# testing two command line parameters

total=$(( $1 * $2 ))
echo The first parameter is $1.
echo The second parameter is $2.
echo The total value is $total.
$
$ ./test2 2 5
The first parameter is 2.
The second parameter is 5.
The total value is 10.
$

```

The shell assigns each parameter to the appropriate variable.

In this example, the command line parameters used were both numerical values. You can also use text strings in the command line:

```

$ cat test3
#!/bin/bash
# testing string parameters

```

```
echo Hello $1, glad to meet you.  
$  
$ ./test3 Rich  
Hello Rich, glad to meet you.  
$
```

The shell passes the string value entered into the command line to the script. However, you'll have a problem if you try to do this with a text string that contains spaces:

```
$ ./test3 Rich Blum  
Hello Rich, glad to meet you.  
$
```

Remember that each of the parameters is separated by a space, so the shell interpreted the space as just separating the two values. To include a space as a parameter value, you must use quotation marks (either single or double quotation marks):

```
$ ./test3 'Rich Blum'  
Hello Rich Blum, glad to meet you.  
$  
$ ./test3 "Rich Blum"  
Hello Rich Blum, glad to meet you.  
$
```

Note

The quotation marks used when you pass text strings as parameters are not part of the data. They just delineate the beginning and the end of the data.

If your script needs more than nine command line parameters, you can continue, but the variable names change slightly. After the ninth variable, you must use braces around the variable number, such as `${10}`. Here's

an example of doing that:

```
$ cat test4
#!/bin/bash
# handling lots of parameters

total=$(( ${10} * ${11} ])
echo The tenth parameter is ${10}
echo The eleventh parameter is ${11}
echo The total is $total
$
$ ./test4 1 2 3 4 5 6 7 8 9 10 11 12
The tenth parameter is 10
The eleventh parameter is 11
The total is 110
$
```

This technique allows you to add as many command line parameters to your scripts as you could possibly need.

Reading the Program Name

You can use the `$0` parameter to determine the name of the program that the shell started from the command line. This can come in handy if you're writing a utility that can have multiple functions. However, there's a small problem that you'll have to deal with. Look what happens in this simple example:

```
$ cat test5
#!/bin/bash
# testing the $0 parameter

echo The command entered is: $0
$
$ ./test5
The command entered is: ./test5
$
$ /home/rich/test5
The command entered is: /home/rich/test5
$
```

When the actual string passed in the `$0` variable is the

entire script path, then the entire path will be used for the program, and not just the program name.

If you want to write a script that performs different functions based on the name of the script run from the command line, you'll have to do a little work. You need to be able to strip off whatever path is used to run the script from the command line.

Fortunately, there's a handy little command available for you that does just that. The `basename` command returns just the program name without the path. Let's modify the example script and see how this works:

```
$ cat test5b
#!/bin/bash
# using basename with the $0 parameter

name='basename $0'
echo The command entered is: $name
$
$ ./test5b
The command entered is: test5b
$
$ /home/rich/test5b
The command entered is: test5b
$
```

Now that's much better. You can now use this technique to write scripts that perform different functions based on the script name used. Here's a simple example to demonstrate this:

```
$ cat test6
#!/bin/bash
# testing a multi-function script

name='basename $0'

if [ $name = "addem" ]
then
    total=$(( $1 + $2 )
```

```

elif [ $name = "multm" ]
then
    total=$(( $1 * $2 ))
fi
echo The calculated value is $total
$
$ chmod u+x test6
$ cp test6 addem
$ ln -s test6 multem
$ ls -l
-rwxr--r--    1 rich      rich      211 Oct 15 18:00 addem
lrwxrwxrwx    1 rich      rich        5 Oct 15 18:01 multem -> test6
-rwxr--r--    1 rich      rich      211 Oct 15 18:00 test6
$
$ ./addem 2 5
The calculated value is 7
$
$ ./multem 2 5
The calculated value is 10
$

```

The example creates two separate filenames from the `test6` code, one by just copying the file and the other by using a link to create the new file. In both cases, the script determines the base name of the script and performs the appropriate function based on that value.

Testing Parameters

You need to be careful when using command line parameters in your shell scripts. If the script runs without the parameters, bad things can happen:

```

$ ./addem 2
./addem: line 8: 2 + : syntax error: operand expected (error
token is " ")
The calculated value is
$

```

When the script assumes there is data in a parameter variable, and there isn't, most likely you'll get an error

message from your script. This is a poor way to write scripts. It's always a good idea to check your parameters to make sure the data is there before using it:

```
$ cat test7
#!/bin/bash
# testing parameters before use

if [ -n "$1" ]
then
    echo Hello $1, glad to meet you.
else
    echo "Sorry, you did not identify yourself. "
fi
$
$ ./test7 Rich
Hello Rich, glad to meet you.
$
$ ./test7
Sorry, you did not identify yourself.
$
```

In this example, the `-n` parameter was used in the `test` command to check if there was data in the command line parameter. In the next section, you'll see there is yet another way to check for command line parameters.

Special Parameter Variables

There are a few special variables available in the bash shell, which track command line parameters. This section describes what they are, and how to use them.

Counting Parameters

As you saw in the last section, it's often a good idea to verify command line parameters before using them in your

script. For scripts that use multiple command line parameters, this can get tedious.

Instead of testing each parameter, you can just count how many parameters were entered on the command line. The bash shell provides a special variable for this purpose.

The special `$#` variable contains the number of command line parameters included when the script was run. You can use this special variable anywhere in the script, just like a normal variable:

```
$ cat test8
#!/bin/bash
# getting the number of parameters

echo There were $# parameters supplied.
$
$ ./test8
There were 0 parameters supplied.
$
$ ./test8 1 2 3 4 5
There were 5 parameters supplied.
$
$ ./test8 1 2 3 4 5 6 7 8 9 10
There were 10 parameters supplied.
$
$ ./test8 "Rich Blum"
There were 1 parameters supplied.
$
```

Now you have the ability to test the number of parameters present before trying to use them:

```
$ cat test9
#!/bin/bash
# testing parameters

if [ $# -ne 2 ]
then
    echo Usage: test9 a b
else
    total=$(( $1 + $2 ))
    echo The total is $total
fi
```

```
$ ./test9
Usage: test9 a b
$
$ ./test9 10
Usage: test9 a b
$
$ ./test9 10 15
The total is 25
$
$ ./test9 10 15 20
Usage: test9 a b
$
```

The `if-then` statement uses the `test` command to perform a numeric test of the number of parameters supplied on the command line. If the correct number of parameters isn't present, you can print an error message that shows the correct usage of the script.

This variable also provides a cool way of grabbing the last parameter on the command line, without having to know how many parameters were used. However, you need to use a little trick to get there.

If you think this through, you might think that because the `$#` variable contains the value of the number of parameters, then using the variable `${ $# }` would represent the last command line parameter variable. Try that out and see what happens:

```
$ cat badtest1
#!/bin/bash
# testing grabbing last parameter

echo The last parameter was ${ $# }
$
$ ./badtest1 10
The last parameter was 15354
$
```

Wow, what happened here? Obviously, something wrong

happened. It turns out that you can't use the dollar sign within the braces. Instead, you must replace the dollar sign with an exclamation mark. Odd, but it works:

```
$ cat test10
#!/bin/bash
# grabbing the last parameter

params=$#
echo The last parameter is $params
echo The last parameter is ${!#}
$
$ ./test10 1 2 3 4 5
The last parameter is 5
The last parameter is 5
$
$ ./test10
The last parameter is 0
The last parameter is ./test10
$
```

Perfect. This test also assigned the `$#` variable value to the variable *params* and then used that variable within the special command line parameter variable format as well. Both versions worked. It's also important to notice that, when there weren't any parameters on the command line, the `$#` value was zero, which is what appears in the *params* variable, but the `${!#}` variable returns the script name used on the command line.

Grabbing All the Data

There are situations where you'll want to just grab all of the parameters provided on the command line and iterate through all of them. Instead of having to mess with using the `$#` variable to determine how many parameters are on the command line, then having to loop through all of them, you

can use a couple of other special variables.

The `$*` and `$@` variables provide easy access to all of your parameters. Both of these variables include all of the command line parameters within a single variable.

The `$*` variable takes all of the parameters supplied on the command line as a single word. The word contains each of the values as they appear on the command line. Basically, instead of treating the parameters as multiple objects, the `$*` variable treats them all as one parameter.

The `$@` variable, on the other hand, takes all of the parameters supplied on the command line as separate words in the same string. It allows you to iterate through the value, separating out each parameter supplied. This is most often accomplished using the `for` command.

It can easily get confusing to figure out how these two variables operate. Let's take a look, so you can see the difference between the two:

```
$ cat testll
#!/bin/bash
# testing $* and $@

echo "Using the \$* method: $*"
echo "Using the \$@ method: $@"
$
$ ./testll rich barbara katie jessica
Using the $* method: rich barbara katie jessica
Using the $@ method: rich barbara katie jessica
$
```

Notice that on the surface, both variables produce the same output, showing all of the command line parameters provided at once.

The following example demonstrates where the

differences are:

```
$ cat test12
#!/bin/bash
# testing $* and $@

count=1
for param in "$*"
do
    echo "\$* Parameter #$count = $param"
    count=$(( $count + 1 ))
done

count=1
for param in "$@"
do
    echo "\$@ Parameter #$count = $param"
    count=$(( $count + 1 ))
done
$
$ ./test12 rich barbara katie jessica
$* Parameter #1 = rich barbara katie jessica
$@ Parameter #1 = rich
$@ Parameter #2 = barbara
$@ Parameter #3 = katie
$@ Parameter #4 = jessica
$
```

Now we're getting somewhere. By using the `for` command to iterate through the special variables, you can see how they each treat the command line parameters differently. The `$*` variable treated all of the parameters as a single parameter, while the `$@` variable treated each parameter separately. This is a great way to iterate through command line parameters.

Being Shifty

Another tool you have in your bash shell tool belt is the `shift` command. The bash shell provides the `shift` command to

help you manipulate command line parameters. The `shift` command literally shifts the command line parameters in their relative positions.

When you use the `shift` command, it “downgrades” each parameter variable one position by default. Thus, the value for variable `$3` is moved to `$2`, the value for variable `$2` is moved to `$1`, and the value for variable `$1` is discarded (note that the value for variable `$0`, the program name, remains unchanged).

This is another great way to iterate through command line parameters, especially if you don't know how many parameters are available. You can just operate on the first parameter, shift the parameters over, and then operate on the first parameter again.

Here's a short demonstration of how this works:

```
$ cat test13
#!/bin/bash
# demonstrating the shift command

count=1
while [ -n "$1" ]
do
    echo "Parameter #$count = $1"
    count=$(( $count + 1 ))
    shift
done

$
$ ./test13 rich barbara katie jessica
Parameter #1 = rich
Parameter #2 = barbara
Parameter #3 = katie
Parameter #4 = jessica
$
```

The script performs a `while` loop, testing the length of the first parameter's value. When the first parameter's length is

zero, the loop ends.

After testing the first parameter, the `shift` command is used to shift all of the parameters one position.

Alternatively, you can perform a multiple location shift by providing a parameter to the `shift` command. Just provide the number of places you want to shift:

```
$ cat test14
#!/bin/bash
# demonstrating a multi-position shift

echo "The original parameters: $*"
shift 2
echo "Here's the new first parameter: $1"
$
$ ./test14 1 2 3 4 5
The original parameters: 1 2 3 4 5
Here's the new first parameter: 3
$
```

By using values in the `shift` command, you can easily skip over parameters you don't need.

Caution

Be careful when working with the `shift` command. When a parameter is shifted out, its value is lost and can't be recovered.

Working with Options

If you've been following along in the book, you've seen several bash commands that provide both parameters and options. *Options* are single letters preceded by a dash that

alter the behavior of a command. This section shows three different methods for working with options in your shell scripts.

Finding Your Options

On the surface, there's nothing all that special about command line options. They appear on the command line immediately after the script name, just the same as command line parameters. In fact, if you want, you can process command line options the same way that you process command line parameters.

Processing Simple Options

In the `test13` script earlier, you saw how to use the `shift` command to work your way down the command line parameters provided with the script program. You can use this same technique to process command line options.

As you extract each individual parameter, use the `case` statement to determine when a parameter is formatted as an option:

```
$ cat test15
#!/bin/bash
# extracting command line options as parameters

while [ -n "$1" ]
do
    case "$1" in
        -a) echo "Found the -a option" ;;
        -b) echo "Found the -b option" ;;
        -c) echo "Found the -c option" ;;
        *) echo "$1 is not an option" ;;
    esac
done
```

```
shift
done
$
$ ./test15 -a -b -c -d
Found the -a option
Found the -b option
Found the -c option
-d is not an option
$
```

The `case` statement checks each parameter for valid options. When one is found, the appropriate commands are run in the `case` statement.

This method works, no matter what order the options are presented on the command line:

```
$ ./test15 -d -c -a
-d is not an option
Found the -c option
Found the -a option
$
```

The `case` statement processes each option as it finds it in the command line parameters. If any other parameters are included on the command line, you can include commands in the catch-all part of the `case` statement to process them.

Separating Options from Parameters

Often you'll run into situations where you'll want to use both options and parameters for a shell script. The standard way to do this in Linux is to separate the two with a special character code that tells the script when the options are done and when the normal parameters start.

For Linux, this special character is the double dash (`--`). The shell uses the double dash to indicate the end of the option list. After seeing the double dash, your script can

safely process the remaining command line parameters as parameters and not options.

To check for the double dash, simply add another entry in the `case` statement:

```
$ cat test16
#!/bin/bash
# extracting options and parameters

while [ -n "$1" ]
do
    case "$1" in
        -a) echo "Found the -a option" ;;
        -b) echo "Found the -b option";;
        -c) echo "Found the -c option" ;;
        --) shift
            break ;;
        *) echo "$1 is not an option";;
    esac
    shift
done

count=1
for param in @$
do
    echo "Parameter #$count: $param"
    count=$((count + 1))
done
$
```

This script uses the `break` command to break out of the `while` loop when it encounters the double dash. Because we're breaking out prematurely, we need to ensure that we stick in another `shift` command to get the double dash out of the parameter variables.

For the first test, try running the script using a normal set of options and parameters:

```
$ ./test16 -c -a -b test1 test2 test3
Found the -c option
Found the -a option
Found the -b option
test1 is not an option
```

```
test2 is not an option
test3 is not an option
$
```

The results show that the script assumed that all the command line parameters were options when it processed them. Next, try the same thing, only this time using the double dash to separate the options from the parameters on the command line:

```
$ ./test16 -c -a -b -- test1 test2 test3
Found the -c option
Found the -a option
Found the -b option
Parameter #1: test1
Parameter #2: test2
Parameter #3: test3
$
```

When the script reaches the double dash, it stops processing options and assumes that any remaining parameters are command line parameters.

Processing Options with Values

Some options require an additional parameter value. In these situations, the command line looks something like this:

```
$ ./testing -a test1 -b -c -d test2
```

Your script must be able to detect when your command line option requires an additional parameter and be able to process it appropriately. Here's an example of how to do that:

```
$ cat test17
#!/bin/bash
# extracting command line options and values
```

```

while [ -n "$1" ]
do
    case "$1" in
        -a) echo "Found the -a option";;
        -b) param="$2"
            echo "Found the -b option, with parameter value $param"
            shift 2;;
        -c) echo "Found the -c option";;
        --) shift
            break;;
        *) echo "$1 is not an option";;
    esac
    shift
done

count=1
for param in "$@"
do
    echo "Parameter #$count: $param"
    count=$(( $count + 1 ))
done
$
$ ./test17 -a -b test1 -d
Found the -a option
Found the -b option, with parameter value test1
-d is not an option
$

```

In this example, the `case` statement defines three options that it processes. The `-b` option also requires an additional parameter value. Since the parameter being processed is `$1`, you know that the additional parameter value is located in `$2` (because all of the parameters are shifted after they are processed). Just extract the parameter value from the `$2` variable. Of course, because we used two parameter spots for this option, you also need to set the `shift` command to shift two positions.

Just as with the basic feature, this process works no matter what order you place the options in (just remember to include the appropriate option parameter with the each option):


```
$ ./test17 -b test1 -a -d
Found the -b option, with parameter value test1
Found the -a option
-d is not an option
$
```

Now you have the basic ability to process command line options in your shell scripts, but there are limitations. For example, this won't work if you try to combine multiple options in one parameter:

```
$ ./test17 -ac
-ac is not an option
$
```

It is a common practice in Linux to combine options, and if your script is going to be user-friendly, you'll want to offer this feature for your users as well. Fortunately, there's another method for processing options that can help you.

Using the getopt Command

The `getopt` command is a great tool to have handy when processing command line options and parameters. It reorganizes the command line parameters to make parsing them in your script easier.

The Command Format

The `getopt` command can take a list of command line options and parameters, in any form, and automatically turn them into the proper format. It uses the following command format:

```
getopt options optstring parameters
```

The *optstring* is the key to the process. It defines the valid option letters used in the command line. It also defines which option letters require a parameter value.

First, list each command line option letter you're going to use in your script in the *optstring*. Then, place a colon after each option letter that requires a parameter value. The `getopt` command parses the supplied parameters based on the *optstring* you define.

Here's a simple example of how `getopt` works:

```
$ getopt ab:cd -a -b test1 -cd test2 test3
-a -b test1 -c -d -- test2 test3
$
```

The *optstring* defines four valid option letters, a, b, c, and d. It also defines that the option letter b requires a parameter value. When the `getopt` command runs, it examines the provided parameter list and parses it based on the supplied *optstring*. Notice that it automatically separated the `-cd` options into two separate options and inserted the double dash to separate the additional parameters on the line.

If you specify an option not in the *optstring*, by default the `getopt` command produces an error message:

```
$ getopt ab:cd -a -b test1 -cde test2 test3
getopt: invalid option -- e
-a -b test1 -c -d -- test2 test3
$
```

If you prefer to just ignore the error messages, use the `-q` option with the command:

```
$ getopt -q ab:cd -a -b test1 -cde test2 test3
-a -b 'test1' -c -d -- 'test2' 'test3'
$
```

Note that the `getopt` command options must be listed before the `optstring`. Now you should be ready to use this command in your scripts to process command line options.

Using getopt in Your Scripts

You can use the `getopt` command in your scripts to format any command line options or parameters entered for your script. It's a little tricky, however, to use.

The trick is to replace the existing command line options and parameters with the formatted version produced by the `getopt` command. The way to do that is to use the `set` command.

You saw the `set` command back in Chapter 5. The `set` command works with the different variables in the shell. Chapter 5 showed how to use the `set` command to display all of the system environment variables.

One of the options of the `set` command is the double dash, which instructs it to replace the command line parameter variables with the values on the `set` command's command line.

The trick then is to feed the original script command line parameters to the `getopt` command, and then feed the output of the `getopt` command to the `set` command to replace the original command line parameters with the nicely formatted ones from `getopt`. This looks something like this:

```
set -- `getopts -q ab:cd "$@"`
```

Now the values of the original command line parameter

variables are replaced with the output from the `getopt` command, which formats the command line parameters for us.

Using this technique, we can now write scripts that handle our command line parameters for us:

```
$ cat test18
#!/bin/bash
# extracting command line options and values with getopt

set -- `getopt -q ab:c "$@"`
while [ -n "$1" ]
do
    case "$1" in
        -a) echo "Found the -a option" ;;
        -b) param="$2"
            echo "Found the -b option, with parameter value $param"
            shift ;;
        -c) echo "Found the -c option" ;;
        --) shift
            break;;
        *) echo "$1 is not an option";;
    esac
    shift
done

count=1
for param in "$@"
do
    echo "Parameter #$count: $param"
    count=$(( $count + 1 ))
done
$
```

You'll notice this is basically the same script as in `test17`. The only thing that changed is the addition of the `getopt` command to help format our command line parameters.

Now when you run the script with complex options, things work much better:

```
$ ./test18 -ac
Found the -a option
Found the -c option
$
```

And of course, all of the original features work just fine as well:

```
$ ./test18 -a -b test1 -cd test2 test3 test4
Found the -a option
Found the -b option, with parameter value 'test1'
Found the -c option
Parameter #1: 'test2'
Parameter #2: 'test3'
Parameter #3: 'test4'
$
```

Now things are looking pretty fancy. However, there's still one small bug that lurks in the `getopt` command. Check out this example:

```
$ ./test18 -a -b test1 -cd "test2 test3" test4
Found the -a option
Found the -b option, with parameter value 'test1'
Found the -c option
Parameter #1: 'test2'
Parameter #2: 'test3'
Parameter #3: 'test4'
$
```

The `getopt` command isn't good at dealing with parameter values with spaces. It interpreted the space as the parameter separator, instead of following the double quotation marks and combining the two values into one parameter. Fortunately, there's yet another solution that solves this problem.

The More Advanced `getopts`

The `getopts` command (notice that it is plural) is built into the bash shell. It looks a lot like its `getopt` cousin, but has some expanded features.

Unlike `getopt`, which produces one output for all of the

processed options and parameters found in the command line, the `getopts` command works on the existing shell parameter variables sequentially.

It processes the parameters it detects in the command line one at a time each time it's called. When it runs out of parameters, it exits with an exit status greater than zero. This makes it great for using in loops to parse all of the parameters on the command line.

The format of the `getopts` command is:

```
getopts optstring variable
```

The `optstring` value is similar to the one used in the `getopt` command. Valid option letters are listed in the `optstring`, along with a colon if the option letter requires a parameter value. To suppress error messages, start the `optstring` with a colon. The `getopts` command places the current parameter in the `variable` defined in the command line.

There are two environment variables that the `getopts` command uses. The `OPTARG` environment variable contains the value to be used if an option requires a parameter value. The `OPTIND` environment variable contains the value of the current location within the parameter list where `getopts` left off. This allows you to continue processing other command line parameters after finishing the options.

Let's take a look at a simple example that uses the `getopts` command:

```
$ cat test19
#!/bin/bash
# simple demonstration of the getopts command

while getopts :ab:c opt
```

```
done
case "$opt" in
a) echo "Found the -a option" ;;
b) echo "Found the -b option, with value $OPTARG";;
c) echo "Found the -c option" ;;
*) echo "Unknown option: $opt";;
esac

$
$ ./test19 -ab test1 -c
Found the -a option
Found the -b option, with value test1
Found the -c option
$
```

The `while` statement defines the `getopts` command, specifying what command line options to look for, along with the variable name to store them in for each iteration.

You'll notice something different about the `case` statement in this example. When the `getopts` command parses the command line options, it also strips off the leading dash, so you don't need them in the `case` definitions.

There are several nice features in the `getopts` command. For starters, you can now include spaces in your parameter values:

```
$ ./test19 -b "test1 test2" -a
Found the -b option, with value test1 test2
Found the -a option
$
```

Another nice feature is that you can run the option letter and the parameter value together without a space:

```
$ ./test19 -abtest1
Found the -a option
Found the -b option, with value test1
$
```

The `getopts` command correctly parsed the `test1` value from the `-b` option. Yet another nice feature of the `getopts`

command is that it bundles any undefined option that it finds in the command line into a single output, the question mark:

```
$ ./test19 -d
Unknown option: ?
$
$ ./test19 -acde
Found the -a option
Found the -c option
Unknown option: ?
Unknown option: ?
$
```

Any option letter not defined in the *optstring* value is sent to your code as a question mark.

The `getopts` command knows when to stop processing options, and leave the parameters for you to process. As `getopts` processes each option, it increments the *OPTIND* environment variable by one. When you've reached the end of the `getopts` processing, you can just use the *OPTIND* value with the `shift` command to move to the parameters:

```
$ cat test20
#!/bin/bash
# processing options and parameters with getopts

while getopts :ab:cd opt
do
    case "$opt" in
        a) echo "Found the -a option" ;;
        b) echo "Found the -b option, with value $OPTARG";;
        c) echo "Found the -c option";;
        d) echo "Found the -d option";;
        *) echo "Unknown option: $opt";;
    esac
done
shift $[ $OPTIND - 1 ]

count=1
for param in "$@"
do
    echo "Parameter $count: $param"
```



```
count=$((count + 1))
done
$
$ ./test20 -a -b test1 -d test2 test3 test4
Found the -a option
Found the -b option, with value test1
Found the -d option
Parameter 1: test2
Parameter 2: test3
Parameter 3: test4
$
```

Now you have a full-featured command line option and parameter processing utility you can use in all of your shell scripts.

Standardizing Options

When you create your shell script, obviously you're in control of what happens. It's completely up to you as to which letter options you select to use and how you select to use them.

However, there are a few letter options that have achieved somewhat of a standard meaning in the Linux world. If you leverage these options in your shell script, it will make your scripts more user-friendly.

[Table 13.1](#) shows some of the common meanings for command line options used in Linux.

[Table 13.1](#) Common Linux Command Line Options

Option	Description
-a	Show all objects.
-c	Produce a count.

-d	Specify a directory.
-e	Expand an object.
-f	Specify a file to read data from.
-h	Display a help message for the command.
-i	Ignore text case.
-l	Produce a long format version of the output.
-n	Use a non-interactive (batch) mode.
-o	Specify an output file to redirect all output to.
-q	Run in quiet mode.
-r	Process directories and files recursively.
-s	Run in silent mode.
-v	Produce verbose output.
-x	Exclude and object.
-y	Answer yes to all questions.

You'll probably recognize most of these option meanings just from working with the various bash commands throughout the book. Using the same meaning for your options helps users interact with your script without having to worry about manuals.

Getting User Input

While providing command line options and parameters is a great way to get data from your script users, sometimes your script needs to be more interactive. There are times when you need to ask a question while the script is running

and wait for a response from the person running your script. The `bash` shell provides the `read` command just for this purpose.

Basic Reading

The `read` command accepts input from the standard input (the keyboard) or from another file descriptor (see Chapter 14). After receiving the input, the `read` command places the data into a standard variable. Here's the `read` command at its simplest:

```
$ cat test21
#!/bin/bash
# testing the read command

echo -n "Enter your name: "
read name
echo "Hello $name, welcome to my program. "
$
$ ./test21
Enter your name: Rich Blum
Hello Rich Blum, welcome to my program.
$
```

That's pretty simple. Notice that the `echo` command that produced the prompt uses the `-n` option. This suppresses the newline character at the end of the string, allowing the script user to enter data immediately after the string, instead of on the next line. This gives your scripts a more form-like appearance.

In fact, the `read` command includes the `-p` option, which allows you to specify a prompt directly in the `read` command line:

```
$ cat test22
#!/bin/bash
```

```
testing the read -p option
```

```
read -p "Please enter your age: " age
days=$(( $age * 365 ))
echo "That makes you over $days days old! "
$
$ ./test22
Please enter your age:10
That makes you over 3650 days old!
$
```

You'll notice in the first example that when a name was entered, the `read` command assigned both the first name and last name to the same variable. The `read` command will assign all data entered at the prompt to a single variable, or you can specify multiple variables. Each data value entered is assigned to the next variable in the list. If the list of variables runs out before the data does, the remaining data is assigned to the last variable:

```
$ cat test23
#!/bin/bash
# entering multiple variables

read -p "Enter your name: " first last
echo "Checking data for $last, $first..."
$
$ ./test23
Enter your name: Rich Blum
Checking data for Blum, Rich...
$
```

You can also specify no variables on the `read` command line. If you do that, the `read` command places any data it receives in the special environment variable `REPLY`:

```
$ cat test24
#!/bin/bash
# testing the REPLY environment variable

read -p "Enter a number: "
factorial=1
for (( count=1; count <= $REPLY; count++ ))
do
    factorial=$(( $factorial * $count ))
done
```

```
done
echo "The factorial of $REPLY is $factorial"
$
$ ./test24
Enter a number: 5
The factorial of 5 is 120
$
```

The `REPLY` environment variable will contain all of the data entered in the input, and it can be used in the shell script as any other variable.

Timing Out

There's a danger when using the `read` command. It's quite possible that your script will get stuck waiting for the script user to enter data. If the script must go on regardless of whether there was any data entered, you can use the `-t` option to specify a timer. The `-t` option specifies the number of seconds for the `read` command to wait for input. When the timer expires, the `read` command returns a non-zero exit status:

```
$ cat test25
#!/bin/bash
# timing the data entry

if read -t 5 -p "Please enter your name: " name
then
    echo "Hello $name, welcome to my script"
else
    echo
    echo "Sorry, too slow! "
fi
$
$ ./test25
Please enter your name: Rich
Hello Rich, welcome to my script
$
$ ./test25
Please enter your name:
```

```
Sorry, too slow!
```

```
$
```

Since the `read` command exits with a non-zero exit status if the timer expires, it's easy to use the standard structured statements, such as an `if-then` statement or a `while` loop to track what happened. In this example, when the timer expires, the `if` statement fails, and the shell executes the commands in the `else` section.

Instead of timing the input, you can also set the `read` command to count the input characters. When a preset number of characters has been entered, it automatically exits, assigning the entered data to the variable:

```
$ cat test26
#!/bin/bash
# getting just one character of input

read -n1 -p "Do you want to continue [Y/N]? " answer
case $answer in
Y | y) echo
        echo "fine, continue on...";;
N | n) echo
        echo OK, goodbye
        exit;;
esac
echo "This is the end of the script"
$
$ ./test26
Do you want to continue [Y/N]? Y
fine, continue on...
This is the end of the script
$
$ ./test26
Do you want to continue [Y/N]? n
OK, goodbye
$
```

This example uses the `-n` option with the value of 1, instructing the `read` command to accept only a single character before exiting. As soon as you press the single

character to answer, the `read` command accepts the input and passes it to the variable. There's no need to press the Enter key.

Silent Reading

There are times when you need input from the script user, but you don't want that input to display on the monitor. The classic example is when entering passwords, but there are plenty of other types of data that you will need to hide.

The `-s` option prevents the data entered in the `read` command from being displayed on the monitor (actually, the data is displayed, but the `read` command sets the text color to the same as the background color). Here's an example of using the `-s` option in a script:

```
$ cat test27
#!/bin/bash
# hiding input data from the monitor

read -s -p "Enter your password: " pass
echo
echo "Is your password really $pass? "
$
$ ./test27
Enter your password:
Is your password really T3st1ng?
$
```

The data typed at the input prompt doesn't appear on the monitor but is assigned to the variable for use in the script.

Reading from a File

Finally, you can also use the `read` command to read data

stored in a file on the Linux system. Each call to the `read` command reads a single line of text from the file. When there are no more lines left in the file, the `read` command will exit with a non-zero exit status.

The tricky part of this is getting the data from the file to the `read` command. The most common method for doing this is to `pipe` the result of the `cat` command of the file directly to a `while` command that contains the `read` command. Here's an example of how to do this:

```
$ cat test28
#!/bin/bash
# reading data from a file

count=1
cat test | while read line
do
    echo "Line $count: $line"
    count=$(( $count + 1 ))
done
echo "Finished processing the file"
$
$ cat test
The quick brown dog jumps over the lazy fox.
This is a test, this is only a test.
O Romeo, Romeo! Wherefore art thou Romeo?
$
$ ./test28
Line 1: The quick brown dog jumps over the lazy fox.
Line 2: This is a test, this is only a test.
Line 3: O Romeo, Romeo! Wherefore art thou Romeo?
Finished processing the file
$
```

The `while` command loop continues processing lines of the file with the `read` command, until the `read` command exits with a non-zero exit status.

Summary

This chapter showed three different methods for retrieving data from the script user. Command line parameters allow users to enter data directly on the command line when they run the script. The script uses positional parameters to retrieve the command line parameters and assign them to variables.

The `shift` command allows you to manipulate the command line parameters by rotating them within the positional parameters. This command allows you to easily iterate through the parameters without knowing how many parameters are available.

There are three special variables that you can use when working with command line parameters. The shell sets the `$#` variable to the number of parameters entered on the command line. The `$*` variable contains all of the parameters as a single string, and the `$@` variable contains all of the parameters as separate words. These variables come in handy when trying to process long parameter lists.

Besides parameters, your script users can also use command line options to pass information to your script. Command line options are single letters preceded by a dash. Different options can be assigned to alter the behavior of your script. The bash shell provides three ways to handle command line options.

The first way is to handle them just like command line parameters. You can iterate through the options using the

positional parameter variables, processing each option as it appears on the command line.

Another way to handle command line options is with the `getopt` command. This command converts command line options and parameters into a standard format that you can process in your script. The `getopt` command allows you to specify which letters it recognizes as options and which options require an additional parameter value. The `getopt` command processes the standard command line parameters and outputs the options and parameters in the proper order.

The final method for handling command line options is via the `getopts` command (note that it's plural). The `getopts` command provides more advanced processing of the command line parameters. It allows for multi-value parameters, along with identifying options not defined by the script.

An interactive method to obtain data from your script users is the `read` command. The `read` command allows your scripts to query users for information and wait. The `read` command places any data entered by the script user into one or more variables, which you can use within the script.

Several options are available for the `read` command that allow you to customize the data input into your script, such as using hidden data entry, applying timed data entry, and requesting a specific number of input characters.

In the next chapter, we look further into how bash shell scripts output data. So far, you've seen how to display data

on the monitor and redirect it to a file. Next, we explore a few other options that you have available not only to direct data to specific locations but also to direct specific types of data to specific locations. This will help make your shell scripts look professional!

Chapter 14

Presenting Data

In This Chapter

- Revisiting redirection
- Standard input and output
- Reporting errors
- Throwing away data
- Creating log files

So far the scripts shown in this book display information either by echoing data to the monitor or by redirecting data to a file. Chapter 10 demonstrated how to redirect the output of a command to a file. This chapter expands on that topic by showing you how you can redirect the output of your script to different locations on your Linux system.

Understanding Input and Output

So far, you've seen two methods for displaying the output

from your scripts:

- Displaying output on the monitor screen
- Redirecting output to a file

Both methods produced an all-or-nothing approach to data output. There are times, however, when it would be nice to display some data on the monitor and other data in a file. For these instances, it comes in handy to know how Linux handles input and output so that you can get your script output to the right place.

The following sections describe how to use the standard Linux input and output system to your advantage, to help direct script output to specific locations.

Standard File Descriptors

The Linux system handles every object as a file. This includes the input and output process. Linux identifies each file object using a *file descriptor*. The file descriptor is a non-negative integer, which uniquely identifies open files in a session. Each process is allowed to have up to nine open file descriptors at a time. The bash shell reserves the first three file descriptors (0, 1, and 2) for special purposes. These are shown in [Table 14.1](#).

Table 14.1 Linux Standard File Descriptors

File Descriptor	Abbreviation	Description
0	STDIN	Standard input
1	STDOUT	Standard output

2	STDERR	Standard error
---	--------	----------------

These three special file descriptors handle the input and output from your script. The shell uses them to direct the default input and output in the shell to the appropriate location (which by default is usually your monitor). The following sections describe each of these standard file descriptors in greater detail.

STDIN

The `STDIN` file descriptor references the standard input to the shell. For a terminal interface, the standard input is the keyboard. The shell receives input from the keyboard on the `STDIN` file descriptor and processes each character as you type it.

When you use the input redirect symbol (`<`), Linux replaces the standard input file descriptor with the file referenced by the redirection. It reads the file and retrieves data just as if it were typed on the keyboard.

Many bash commands accept input from `STDIN`, especially if no files are specified on the command line. Here's an example of using the `cat` command with data entered from `STDIN`:

```
$ cat
this is a test
this is a test
this is a second test.
this is a second test.
```

When you enter the `cat` command on the command line by itself, it accepts input from `STDIN`. As you enter each line,

the `cat` command echoes the line to the display.

However, you can also use the `STDIN` redirect symbol to force the `cat` command to accept input from another file other than `STDIN`:

```
$ cat < testfile
This is the first line.
This is the second line.
This is the third line.
$
```

Now the `cat` command uses the lines that are contained in the `testfile` file as the input. You can use this technique to input data to any shell command that accepts data from `STDIN`.

STDOUT

The `STDOUT` file descriptor references the standard output for the shell. On a terminal interface, the standard output is the terminal monitor. All output from the shell (including programs and scripts you run in the shell) is directed to the standard output, which is the monitor.

Most bash commands direct their output to the `STDOUT` file descriptor by default. As shown in Chapter 10, you can change that using output redirection:

```
$ ls -l > test2
$ cat test2
total 20
-rw-rw-r-- 1 rich rich 53 2010-10-16 11:30 test
-rw-rw-r-- 1 rich rich  0 2010-10-16 11:32 test2
-rw-rw-r-- 1 rich rich 73 2010-10-16 11:23 testfile
$
```

With the output redirection symbol, all of the output that normally would have gone to the monitor is instead

redirected to the designated redirection file by the shell.

You can also append data to a file. You do this using the `>>` symbol:

```
$ who >> test2
$ cat test2
total 20
-rw-rw-r-- 1 rich rich 53 2010-10-16 11:30 test
-rw-rw-r-- 1 rich rich 0 2010-10-16 11:32 test2
-rw-rw-r-- 1 rich rich 73 2010-10-16 11:23 testfile
rich pts/0 2010-10-17 15:34 (192.168.1.2)
$
```

The output generated by the `who` command is appended to the data already in the `test2` file.

However, if you use the standard output redirection for your scripts, you can run into a problem. Here's an example of what can happen in your script:

```
$ ls -al badfile > test3
ls: cannot access badfile: No such file or directory
$ cat test3
$
```

When a command produces an error message, the shell doesn't redirect the error message to the output redirection file. The shell created the output redirection file, but the error message appeared on the monitor screen. Notice that there isn't an error when trying to display the contents of the `test3` file. The `test3` file was created just fine, it's just empty.

The shell handles error messages separately from the normal output. If you're creating a shell script that runs in background mode, often you must rely on the output messages being sent to a log file. Using this technique, if any error messages occur, they won't appear in the log file.

You'll need to do something different.

STDERR

The shell handles error messages using the special `STDERR` file descriptor. The `STDERR` file descriptor references the standard error output for the shell. This is the location where the shell sends error messages generated by the shell or programs and scripts running in the shell.

By default, the `STDERR` file descriptor points to the same place as the `STDOUT` file descriptor (even though they are assigned different file descriptor values). This means that, by default, all error messages go to the monitor display.

However, as you saw in the example, when you redirect `STDOUT` this doesn't automatically redirect `STDERR`. When working with scripts, you'll often want to change that behavior, especially if you're interested in logging error messages to a log file.

Redirecting Errors

You've already seen how to redirect the `STDOUT` data by using the redirection symbol. Redirecting the `STDERR` data isn't much different; you just need to define the `STDERR` file descriptor when you use the redirection symbol. There are a couple of ways to do this.

Redirecting Just Errors

As you saw in [Table 14.1](#), the `STDERR` file descriptor is set to the value 2. You can select to redirect only error messages by placing this file descriptor value immediately before the redirection symbol. The value must appear immediately before the redirection symbol or it won't work:

```
$ ls -al badfile 2> test4
$ cat test4
ls: cannot access badfile: No such file or directory
$
```

Now when you run the command, the error message doesn't appear on the monitor. Instead, the output file contains any error messages that are generated by the command. Using this method, the shell redirects the error messages only, not the normal data. Here's another example of mixing `STDOUT` and `STDERR` messages in the same output:

```
$ ls -al test badtest test2 2> test5
-rw-rw-r-- 1 rich rich 158 2010-10-16 11:32 test2
$ cat test5
ls: cannot access test: No such file or directory
ls: cannot access badtest: No such file or directory
$
```

The normal `STDOUT` output from the `ls` command still goes to the default `STDOUT` file descriptor, which is the monitor. Because the command redirects file descriptor 2 output (`STDERR`) to an output file, the shell sends any error messages generated directly to the specified redirection file.

Redirecting Errors and Data

If you want to redirect both errors and the normal output, you'll need to use two redirection symbols. You need to

precede each with the appropriate file descriptor for the data you want to redirect and then have them point to the appropriate output file for holding the data:

```
$ ls -al test test2 test3 badtest 2> test6 1> test7
$ cat test6
ls: cannot access test: No such file or directory
ls: cannot access badtest: No such file or directory
$ cat test7
-rw-rw-r-- 1 rich rich 158 2010-10-16 11:32 test2
-rw-rw-r-- 1 rich rich   0 2010-10-16 11:33 test3
$
```

The shell redirects the normal output of the `ls` command that would have gone to `STDOUT` to the `test7` file using the `1>` symbol. Any error messages that would have gone to `STDERR` were redirected to the `test6` file using the `2>` symbol.

You can use this technique to separate normal script output from any error messages that occur in the script. This allows you to easily identify errors without having to wade through thousands of lines of normal output data.

Alternatively, if you want, you can redirect both `STDERR` and `STDOUT` output to the same output file. The bash shell provides a special redirection symbol just for this purpose, the `&>` symbol:

```
$ ls -al test test2 test3 badtest &> test7
$ cat test7
ls: cannot access test: No such file or directory
ls: cannot access badtest: No such file or directory
-rw-rw-r-- 1 rich rich 158 2010-10-16 11:32 test2
-rw-rw-r-- 1 rich rich   0 2010-10-16 11:33 test3
$
```

When you use the `&>` symbol, all of the output generated by the command is sent to the same location, both data and errors. You'll notice that one of the error messages is out of order from what you'd expect. The error message for

the `badtest` file (the last file to be listed) appears second in the output file. The bash shell automatically gives error messages a higher priority than the standard output. This allows you to view the error messages together, rather than scattered throughout the output file.

Redirecting Output in Scripts

You can use the `STDOUT` and `STDERR` file descriptors in your scripts to produce output in multiple locations simply by redirecting the appropriate file descriptors. There are two methods for redirecting output in the script:

- Temporarily redirecting each line
- Permanently redirecting all commands in the script

The following sections describe how each of these methods works.

Temporary Redirections

If you want to purposely generate error messages in your script, you can redirect an individual output line to `STDERR`. All you need to do is use the output redirection symbol to redirect the output to the `STDERR` file descriptor. When you redirect to a file descriptor, you must precede the file descriptor number with an ampersand sign (`&`):

```
echo "This is an error message" >&2
```

This line displays the text wherever the `STDERR` file

descriptor for the script is pointing, instead of the normal `STDOUT`. The following is an example of a script that uses this feature:

```
$ cat test8
#!/bin/bash
# testing STDERR messages

echo "This is an error" >&2
echo "This is normal output"
$
```

If you run the script as normal, you won't notice any difference:

```
$ ./test8
This is an error
This is normal output
$
```

Remember that, by default, Linux directs the `STDERR` output to `STDOUT`. However, if you redirect `STDERR` when running the script, any text directed to `STDERR` in the script will be redirected:

```
$ ./test8 2> test9
This is normal output
$ cat test9
This is an error
$
```

Perfect! The text that's displayed using `STDOUT` appears on the monitor, while the echo statement text sent to `STDERR` is redirected to the output file.

This method is great for generating error messages in your scripts. If someone uses your scripts, they can easily redirect the error messages using the `STDERR` file descriptor, as shown.

Permanent Redirections

If you have lots of data that you're redirecting in your script, it can get tedious having to redirect every `echo` statement. Instead, you can tell the shell to redirect a specific file descriptor for the duration of the script by using the `exec` command:

```
$ cat testl0
#!/bin/bash
# redirecting all output to a file
exec 1>testout

echo "This is a test of redirecting all output"
echo "from a script to another file."
echo "without having to redirect every individual line"
$ ./testl0
$ cat testout
This is a test of redirecting all output
from a script to another file.
without having to redirect every individual line
$
```

The `exec` command starts a new shell and redirects the `STDOUT` file descriptor to a file. All output in the script that goes to `STDOUT` is instead redirected to the file.

You can also redirect the `STDOUT` in the middle of a script:

```
$ cat testl1
#!/bin/bash
# redirecting output to different locations

exec 2>testerror

echo "This is the start of the script"
echo "now redirecting all output to another location"

exec 1>testout

echo "This output should go to the testout file"
echo "but this should go to the testerror file" >&2
$
$ ./testl1
This is the start of the script
```

```
now redirecting all output to another location
$ cat testout
This output should go to the testout file
$ cat testerror
but this should go to the testerror file
$
```

The script uses the `exec` command to redirect any output going to `STDERR` to the file `testerror`. Next, the script uses the `echo` statement to display a few lines to `STDOUT`. After that, the `exec` command is used again to redirect `STDOUT` to the `testout` file. Notice that even when `STDOUT` is redirected, you can still specify the output from an `echo` statement to go to `STDERR`, which in this case is still redirected to the `testerror` file.

This feature can come in handy when you want to redirect the output of just parts of a script to an alternative location, such as an error log. There's just one problem you'll run into when using this.

Once you redirect `STDOUT` or `STDERR`, you can't easily redirect them back to their original location. If you need to switch back and forth with your redirection, there's a trick you'll need to learn. The “Creating Your Own Redirection” section later in this chapter discusses what this trick is and how to use it in your shell scripts.

Redirecting Input in Scripts

You can use the same technique used to redirect `STDOUT` and `STDERR` in your scripts to redirect `STDIN` from the keyboard. The `exec` command allows you to redirect `STDIN` to a file on the Linux system:

```
exec 0< testfile
```

This command informs the shell that it should retrieve input from the file `testfile` instead of `STDIN`. This redirection applies anytime the script requests input. Here's an example of this in action:

```
$ cat test12
#!/bin/bash
# redirecting file input

exec 0< testfile
count=1

while read line
do
    echo "Line #$count: $line"
    count=$(( $count + 1 ))
done
$ ./test12
Line #1: This is the first line.
Line #2: This is the second line.
Line #3: This is the third line.
$
```

Chapter 13 showed you how to use the `read` command to read data entered from the keyboard by a user. By redirecting `STDIN` to a file, when the `read` command attempts to read from `STDIN`, it retrieves data from the file instead of the keyboard.

This is an excellent technique to read data in files for processing in your scripts. A common task for Linux system administrators is to read data from log files for processing. This is the easiest way to accomplish that task.

Creating Your Own

Redirection

When you redirect input and output in your script, you're not limited to the three default file descriptors. I mentioned that you could have up to nine open file descriptors in the shell. The other six file descriptors are numbered from three through eight and are available for you to use as either input or output redirection. You can assign any of these file descriptors to a file and then use them in your scripts as well. This section shows you how to use the other file descriptors in your scripts.

Creating Output File Descriptors

You assign a file descriptor for output by using the `exec` command. As with the standard file descriptors, once you assign an alternative file descriptor to a file location, that redirection stays permanent until you reassign it. Here's a simple example of using an alternative file descriptor in a script:

```
$ cat test13
#!/bin/bash
# using an alternative file descriptor

exec 3>test13out

echo "This should display on the monitor"
echo "and this should be stored in the file" >&3
echo "Then this should be back on the monitor"
$ ./test13
This should display on the monitor
Then this should be back on the monitor
$ cat test13out
and this should be stored in the file
$
```

The script uses the `exec` command to redirect file descriptor 3 to an alternative file location. When the script executes the `echo` statements, they display on `STDOUT` as you would expect. However, the `echo` statements that you redirect to file descriptor 3 go to the alternative file. This allows you to keep normal output for the monitor, and redirect special information to files, such as log files.

You can also use the `exec` command to append data to an existing file instead of creating a new file:

```
exec 3>>test13out
```

Now the output will be appended to the `test13out` file instead of creating a new file.

Redirecting File Descriptors

Now comes the trick to help you bring back a redirected file descriptor. You can assign an alternative file descriptor to a standard file descriptor, and vice versa. This means that you can redirect the original location of `STDOUT` to an alternative file descriptor and then redirect that file descriptor back to `STDOUT`. This might sound somewhat complicated, but in practice it's fairly straightforward. This simple example will clear things up for you:

```
$ cat test14
#!/bin/bash
# storing STDOUT, then coming back to it

exec 3>&1
exec 1>test14out

echo "This should store in the output file"
echo "along with this line."
```

```
exec 1>&3

echo "Now things should be back to normal"
$
$ ./test14
Now things should be back to normal
$ cat test14out
This should store in the output file
along with this line.
$
```

This example is a little crazy so let's walk through it piece by piece. First, the script redirects file descriptor 3 to the current location of file descriptor 1, which is `STDOUT`. This means that any output sent to file descriptor 3 will go to the monitor.

The second `exec` command redirects `STDOUT` to a file. The shell will now redirect any output sent to `STDOUT` directly to the output file. However, file descriptor 3 still points to the original location of `STDOUT`, which is the monitor. If you send output data to file descriptor 3 at this point, it'll still go to the monitor, even though `STDOUT` is redirected.

After sending some output to `STDOUT`, which points to a file, the script then redirects `STDOUT` to the current location of file descriptor 3, which is still set to the monitor. This means that now `STDOUT` is pointing to its original location, the monitor.

While this method can get confusing, it's a common way to temporarily redirect output in script files and then set the output back to the normal settings.

Creating Input File Descriptors

You can redirect input file descriptors exactly the same way as output file descriptors. Save the `STDIN` file descriptor location to another file descriptor before redirecting it to a file, and then when you're done reading the file you can restore `STDIN` to its original location:

```
$ cat test15
#!/bin/bash
# redirecting input file descriptors

exec 6<&0

exec 0< testfile

count=1
while read line
do
    echo "Line #$count: $line"
    count=$(( $count + 1 ))
done
exec 0<&6
read -p "Are you done now? " answer
case $answer in
Y|y) echo "Goodbye";;
N|n) echo "Sorry, this is the end.";;
esac
$ ./test15
Line #1: This is the first line.
Line #2: This is the second line.
Line #3: This is the third line.
Are you done now? y
Goodbye
$
```

In this example, file descriptor 6 is used to hold the location for `STDIN`. The script then redirects `STDIN` to a file. All of the input for the `read` command comes from the redirected `STDIN`, which is now the input file.

When all of the lines have been read, the script returns `STDIN` to its original location by redirecting it to file descriptor 6. The script tests to make sure that `STDIN` is back to normal

by using another read command, which this time waits for input from the keyboard.

Creating a Read/Write File Descriptor

As odd as it may seem, you can also open a single file descriptor for both input and output. You can then use the same file descriptor to both read data from a file and write data to the same file.

You need to be especially careful with this method, however. As you read and write data to and from a file, the shell maintains an internal pointer, indicating where it is in the file. Any reading or writing occurs where the file pointer last left off. This can produce some interesting results if you're not careful. Take a look at this example:

```
$ cat test16
#!/bin/bash
# testing input/output file descriptor

exec 3<> testfile
read line <&3
echo "Read: $line"
echo "This is a test line" >&3
$ cat testfile
This is the first line.
This is the second line.
This is the third line.
$ ./test16
Read: This is the first line.
$ cat testfile
This is the first line.
This is a test line
ine.
This is the third line.
$
```

This example uses the `exec` command to assign file descriptor 3 for both input and output sent to and from the

file `testfile`. Next, it uses the `read` command to read the first line in the file, using the assigned file descriptor, and then it displays the read line of data in `STDOUT`. After that, it uses the `echo` statement to write a line of data to the file opened with the same file descriptor.

When you run the script, at first things look just fine. The output shows that the script read the first line in the `testfile` file. However, if you display the contents of the `testfile` file after running the script, you'll see that the data written to the file overwrote the existing data.

When the script writes data to the file, it starts where the file pointer is located. The `read` command reads the first line of data, so it left the file pointer pointing to the first character in the second line of data. When the `echo` statement outputs data to the file, it places the data at the current location of the file pointer, overwriting whatever data was there.

Closing File Descriptors

If you create new input or output file descriptors, the shell automatically closes them when the script exits. There are situations, however, when you need to manually close a file descriptor before the end of the script.

To close a file descriptor, redirect it to the special symbol `&-`. This is how this looks in the script:

```
exec 3>&-
```

This statement closes file descriptor 3, preventing it from

being used any more in the script. Here's an example of what happens when you try to use a closed file descriptor:

```
$ cat badtest
#!/bin/bash
# testing closing file descriptors

exec 3> testl7file

echo "This is a test line of data" >&3

exec 3>&-

echo "This won't work" >&3
$ ./badtest
./badtest: 3: Bad file descriptor
$
```

Once you close the file descriptor, you can't write any data to it in your script or the shell produces an error message.

There's yet another thing to be careful of when closing file descriptors. If you open the same output file later on in your script, the shell replaces the existing file with a new file. This means that if you output any data, it will overwrite the existing file. Consider the following example of this problem:

```
$ cat testl7
#!/bin/bash
# testing closing file descriptors

exec 3> testl7file

echo "This is a test line of data" >&3
exec 3>&-

cat testl7file

exec 3> testl7file
echo "This'll be bad" >&3
$ ./testl7
This is a test line of data
$ cat testl7file
```

After sending a data string to the `test17file` file and closing the file descriptor, the script uses the `cat` command to display the contents of the file. So far, so good. Next, the script reopens the output file and sends another data string to it. When you display the contents of the output file, all you'll see is the second data string. The shell overwrote the original output file.

Listing Open File Descriptors

With only nine file descriptors available to you, you'd think that it wouldn't be too hard keeping things straight. Sometimes, however, it's easy to get lost when trying to keep track of which file descriptor is redirected where. To help you keep your sanity, the bash shell provides the `lsof` command.

The `lsof` command lists all of the open file descriptors on the entire Linux system. This is somewhat of a controversial feature, as it can provide information about the Linux system to non-system-administrators. Because of this, many Linux systems hide this command so that users don't accidentally stumble across it.

On my Fedora Linux system, the `lsof` command is located in the `/usr/sbin` directory. To run it with a normal user account, I have to reference it by its full pathname:

This produces an amazing amount of output. It displays information about every file currently open on the Linux system. This includes all of the processes running on background, as well as any user accounts logged in to the system.

There are plenty of command line parameters and options available to help filter out the `lsuf` output. The most commonly used are `-p`, which allows you to specify a process ID (PID), and `-d`, which allows you to specify the file descriptor numbers to display.

To easily determine the current PID of the process, you can use the special environment variable `$$`, which the shell sets to the current PID. The `-a` option is used to perform a Boolean AND of the results of the other two options, to produce the following:

```
$ /usr/sbin/lsuf -a -p $$ -d 0,1,2
COMMAND  PID USER  FD   TYPE DEVICE SIZE NODE NAME
bash     3344 rich   0u    CHR  136,0      2 /dev/pts/0
bash     3344 rich   1u    CHR  136,0      2 /dev/pts/0
bash     3344 rich   2u    CHR  136,0      2 /dev/pts/0
$
```

This shows the default file descriptors (0, 1, and 2) for the current process (the bash shell). The default output of `lsuf` contains several columns of information, described in [Table 14.2](#).

Table 14.2 Default lsuf Output

Column	Description
COMMAND	The first nine characters of the name of the command in the process
PID	The process ID of the process

USER	The login name of the user who owns the process
FD	The file descriptor number and access type (r—(read), w—(write), u—(read/write))
TYPE	The type of file (CHR—(character), BLK—(block), DIR—(directory), REG—(regular file))
DEVICE	The device numbers (major and minor) of the device
SIZE	If available, the size of the file
NODE	The node number of the local file
NAME	The name of the file

The file type associated with `STDIN`, `STDOUT`, and `STDERR` is character mode. Because the `STDIN`, `STDOUT`, and `STDERR` file descriptors all point to the terminal, the name of the output file is the device name of the terminal. All three standard files are available for both reading and writing (although it does seem odd to be able to write to `STDIN` and read from `STDOUT`).

Now, let's take a look at the results of the `lsuf` command from inside a script that's opened a couple of alternative file descriptors:

```
$ cat testl8
#!/bin/bash
# testing lsuf with file descriptors

exec 3> testl8file1
exec 6> testl8file2
exec 7< testfile

/usr/sbin/lsuf -a -p $$ -d0,1,2,3,6,7
$ ./testl8
COMMAND PID USER  FD  TYPE DEVICE SIZE  NODE NAME
testl8  3594 rich   0u   CHR  136,0        2 /dev/pts/0
testl8  3594 rich   1u   CHR  136,0        2 /dev/pts/0
estl8   3594 rich   2u   CHR  136,0        2 /dev/pts/0
18     3594 rich   3w   REG  253,0        0 360712 /home/rich/testl8file1
```

```
18 3594 rich      6w REG 253,0      0 360715 /home/rich/test18file2
18 3594 rich      7r REG 253,0      73 360717 /home/rich/testfile
$
```

The script creates three alternative file descriptors, two for output (3 and 6) and one for input (7). When the script runs the `ls -l` command, you can see the new file descriptors in the output. We truncated the first part of the output so that you could see the results of the file name. The file name shows the complete pathname for the files used in the file descriptors. It shows each of the files as type `REG`, which indicates that they are regular files on the filesystem.

Suppressing Command Output

There are times when you don't want to display any output from your script. This often occurs if you're running a script as a background process (see Chapter 15). If any error messages occur from the script while it's running in the background, the shell e-mails them to the owner of the process. This can get tedious, especially if you run scripts that generate minor nuisance errors.

To solve that problem, you can redirect `STDERR` to a special file called the *null file*. The null file is pretty much what it says it is, a file that contains nothing. Any data that the shell outputs to the null file is not saved, thus they are lost.

The standard location for the null file on Linux systems is `/dev/null`. Any data you redirect to that location is thrown

away and doesn't appear:

```
$ ls -al > /dev/null
$ cat /dev/null
$
```

This is a common way to suppress any error messages without actually saving them:

```
$ ls -al badfile test16 2> /dev/null
-rwxr--r--  1 rich    rich      135 Oct 29 19:57 test16*
$
```

You can also use the `/dev/null` file for input redirection as an input file. Because the `/dev/null` file contains nothing, it is often used by programmers to quickly remove data from an existing file without having to remove the file and re-create it:

```
$ cat testfile
This is the first line.
This is the second line.
This is the third line.
$ cat /dev/null > testfile
$ cat testfile
$
```

The file `testfile` still exists on the system, but now it is empty. This is a common method used to clear out log files that must remain in place for applications to operate.

Using Temporary Files

The Linux system contains a special directory location reserved for temporary files. Linux uses the `/tmp` directory for files that don't need to be kept indefinitely. Most Linux distributions configure the system to automatically remove

any files in the `/tmp` directory at bootup.

Any user account on the system has privileges to read and write files in the `/tmp` directory. This feature provides an easy way for you to create temporary files that you don't necessarily have to worry about cleaning up.

There's even a specific command to use for creating a temporary file. The `mktemp` command allows you to easily create a unique temporary file in the `/tmp` folder. The shell creates the file but doesn't use your default umask value (see Chapter 6). Instead, it only assigns read and write permissions to the file's owner and makes you the owner of the file. Once you create the file, you have full access to read and write to and from it from your script, but no one else will be able to access it (other than the root user of course).

Creating a Local Temporary File

By default, `mktemp` creates a file in the local directory. To create a temporary file in a local directory with the `mktemp` command, all you need to do is specify a file name template. The template consists of any text filename, plus six X's appended to the end of the file name:

```
$ mktemp testing.XXXXXX
$ ls -al testing*
-rw----- 1 rich      rich      0 Oct 17 21:30 testing.UfIi13
$
```

The `mktemp` command replaces the six X's with a six-character code to ensure the file name is unique in the directory. You can create multiple temporary files and be

assured that each one is unique:

```
$ mktemp testing.XXXXXX
testing.lDRLuV
$ mktemp testing.XXXXXX
testing.lVBtkW
$ mktemp testing.XXXXXX
testing.PgqNKG
$ ls -l testing*
-rw----- 1 rich rich 0 Oct 17 21:57 testing.lDRLuV
-rw----- 1 rich rich 0 Oct 17 21:57 testing.PgqNKG
-rw----- 1 rich rich 0 Oct 17 21:30 testing.UfIi13
-rw----- 1 rich rich 0 Oct 17 21:57 testing.lVBtkW
$
```

As you can see, the output of the `mktemp` command is the name of the file that it creates. When you use the `mktemp` command in a script, you'll want to save that filename in a variable, so you can refer to it later on in the script:

```
$ cat test19
#!/bin/bash
# creating and using a temp file

tempfile='mktemp test19.XXXXXX'

exec 3>$tempfile

echo "This script writes to temp file $tempfile"

echo "This is the first line" >&3
echo "This is the second line." >&3
echo "This is the last line." >&3
exec 3>&-

echo "Done creating temp file. The contents are:"
cat $tempfile
rm -f $tempfile 2> /dev/null
$ ./test19
This script writes to temp file test19.vCHoya
Done creating temp file. The contents are:
This is the first line
This is the second line.
This is the last line.
$ ls -al test19*
-rwxr--r-- 1 rich rich 356 Oct 29 22:03 test19*
$
```

The script uses the `mktemp` command to create a temporary file and assigns the file name to the `$tempfile` variable. It then uses the temporary file as the output redirection file for file descriptor 3. After displaying the temporary file name on `STDOUT`, it writes a few lines to the temporary file, and then it closes the file descriptor. Finally, it displays the contents of the temporary file and then uses the `rm` command to remove it.

Creating a Temporary File in /tmp

The `-t` option forces `mktemp` to create the file in the temporary directory of the system. When you use this feature, the `mktemp` command returns the full pathname used to create the temporary file, not just the file name:

```
$ mktemp -t test.XXXXXX
/tmp/test.xG3374
$ ls -al /tmp/test*
-rw----- 1 rich rich 0 2010-10-29 18:41 /tmp/test.xG3374
$
```

Because the `mktemp` command returns the full pathname, you can then reference the temporary file from any directory on the Linux system, no matter where it places the temporary directory:

```
$ cat test20
#!/bin/bash
# creating a temp file in /tmp

tempfile=`mktemp -t tmp.XXXXXX`

echo "This is a test file." > $tempfile
echo "This is the second line of the test." >> $tempfile

echo "The temp file is located at: $tempfile"
cat $tempfile
```

```

rm -f $tempfile
$ ./test20
The temp file is located at: /tmp/tmp.Ma3390
This is a test file.
This is the second line of the test.
$

```

When `mktemp` creates the temporary file, it returns the full pathname to the environment variable. You can then use that value in any command to reference the temporary file.

Creating a Temporary Directory

The `-d` option tells the `mktemp` command to create a temporary directory instead of a file. You can then use that directory for whatever purposes you need, such as creating additional temporary files:

```

$ cat test21
#!/bin/bash
# using a temporary directory

tempdir='mktemp -d dir.XXXXXX'
cd $tempdir
tempfile1='mktemp temp.XXXXXX'
tempfile2='mktemp temp.XXXXXX'
exec 7> $tempfile1
exec 8> $tempfile2

echo "Sending data to directory $tempdir"
echo "This is a test line of data for $tempfile1" >&7
echo "This is a test line of data for $tempfile2" >&8
$ ./test21
Sending data to directory dir.ouT8S8
$ ls -al
total 72
drwxr-xr-x    3 rich      rich          4096 Oct 17 22:20 ./
drwxr-xr-x    9 rich      rich          4096 Oct 17 09:44 ../
drwx-----   2 rich      rich          4096 Oct 17 22:20 dir.ouT8S8/
-rwxr--r--    1 rich      rich          338 Oct 17 22:20 test21*
$ cd dir.ouT8S8
[dir.ouT8S8]$ ls -al
total 16
drwx-----   2 rich      rich          4096 Oct 17 22:20 ./

```



```
drwxr-xr-x    3 rich    rich          4096 Oct 17 22:20 ../
-rw-----    1 rich    rich          44 Oct 17 22:20 temp.N5F306
-rw-----    1 rich    rich          44 Oct 17 22:20 temp.SQslb7

[dir.ouT8S8]$ cat temp.N5F306
This is a test line of data for temp.N5F306
[dir.ouT8S8]$ cat temp.SQslb7
This is a test line of data for temp.SQslb7
[dir.ouT8S8]$
```

The script creates a directory in the current directory and then it uses the `cd` command to change to that directory before creating two temporary files. The two temporary files are then assigned to file descriptors and used to store output from the script.

Logging Messages

Sometimes it's beneficial to send output both to the monitor and to a file for logging. Instead of having to redirect output twice, you can use the special `tee` command.

The `tee` command is like a T-connector for pipes. It sends data from `STDIN` to two destinations at the same time. One destination is `STDOUT`. The other destination is a file name specified on the `tee` command line:

```
tee filename
```

Because `tee` redirects data from `STDIN`, you can use it with the pipe command to redirect output from any command:

```
$ date | tee testfile
Sun Oct 17 18:56:21 EDT 2010
$ cat testfile
Sun Oct 17 18:56:21 EDT 2010
$
```

The output appears in `STDOUT`, and it is written to the file

specified. Be careful: by default, the `tee` command overwrites the output file on each use:

```
$ who | tee testfile
rich      pts/0          2010-10-17 18:41 (192.168.1.2)
$ cat testfile
rich      pts/0          2010-10-17 18:41 (192.168.1.2)
$
```

If you want to append data to the file, you must use the `-a` option:

```
$ date | tee -a testfile
Sun Oct 17 18:58:05 EDT 2010
$ cat testfile
rich      pts/0          2010-10-17 18:41 (192.168.1.2)
Sun Oct 17 18:58:05 EDT 2010
$
```

Using this technique, you can both save data in files and display the data on the monitor for your users:

```
$ cat test22
#!/bin/bash
# using the tee command for logging

tempfile=test22file

echo "This is the start of the test" | tee $tempfile
echo "This is the second line of the test" | tee -a $tempfile
echo "This is the end of the test" | tee -a $tempfile
$ ./test22
This is the start of the test
This is the second line of the test
This is the end of the test
$ cat test22file
This is the start of the test
This is the second line of the test
This is the end of the test
$
```

Now you can save a permanent copy of your output at the same time that you're displaying it to your users.

Summary

Understanding how the bash shell handles input and output can come in handy when creating your scripts. You can manipulate both how the script receives data and how it displays data, to customize your script for any environment. You can redirect the input of a script from the standard input (`STDIN`) to any file on the system. You can also redirect the output of the script from the standard output (`STDOUT`) to any file on the system.

Besides the `STDOUT`, you can redirect any error messages your script generates by redirecting the `STDERR` output. This is accomplished by redirecting the file descriptor associated with the `STDERR` output, which is file descriptor 2. You can redirect `STDERR` output to the same file as the `STDOUT` output or to a completely separate file. This enables you to separate normal script messages from any error messages generated by the script.

The bash shell allows you to create your own file descriptors for use in your scripts. You can create file descriptors 3 through 9 and assign them to any output file you desire. Once you create a file descriptor, you can redirect the output of any command to it, using the standard redirection symbols.

The bash shell also allows you to redirect input to a file descriptor, providing an easy way to read data contained in a file into your script. You can use the `lsof` command to display the active file descriptors in your shell.

Linux systems provide a special file, called `/dev/null`, to allow you to redirect output that you don't want. The Linux system discards anything redirected to the `/dev/null` file. You can also use this file to produce an empty file by redirecting the contents of the `/dev/null` file to the file.

The `mktemp` command is a handy feature of the bash shell that allows you to easily create temporary files and directories. Simply specify a template for the `mktemp` command, and it creates a unique file each time you call it, based on the file template format. You can also create temporary files and directories in the `/tmp` directory on the Linux system, which is a special location that isn't preserved between system boots.

The `tee` command is a handy way to send output both to the standard output and to a log file. This enables you to display messages from your script on the monitor and store them in a log file at the same time.

In Chapter 15, you'll see how to control and run your scripts. Linux provides several different methods for running scripts other than directly from the command line interface prompt. You'll see how to schedule your scripts to run at a specific time, as well as learn how to pause them while they're running.

Chapter 15

Script Control

In This Chapter

- Handling signals
- Running scripts in background mode
- Running scripts without a console
- Job control
- Being nice
- Running like clockwork
- Start at the beginning
- Filesystem commands

As you start building advanced scripts, you'll probably wonder how to run and control them on your Linux system. So far in this book, the only way we've run scripts is directly from the command line interface in real-time mode. This isn't the only way to run scripts in Linux. Quite a few options are available for running your shell scripts. This chapter examines the different ways to get your scripts started. Also, sometimes a script gets stuck in a loop, and you need to figure out how to get it to stop without turning off

your Linux system. This chapter examines the different ways you can control how and when your shell script runs on your system.

Handling Signals

Linux uses signals to communicate with processes running on the system. Chapter 4 described the different Linux signals and how the Linux system uses these signals to stop, start, and kill processes. You can also use these signals to control the operation of your shell script by programming the script to perform certain commands when it receives specific signals from the Linux system.

Linux Signals Revisited

There are over 30 Linux signals that can be generated by the system and applications. [Table 15.1](#) lists the most common Linux system signals that you'll run across in your Linux programming.

Table 15.1 Linux Signals

Signal	Value	Description
1	SIGHUP	Hang up the process.
2	SIGINT	Interrupt the process.
3	SIGQUIT	Stop the process.
9	SIGKILL	Unconditionally terminate the process.
15	SIGTERM	Terminate the process if possible.

17	SIGSTOP	Unconditionally stop, but don't terminate, the process.
18	SIGTSTP	Stop or pause the process, but don't terminate.
19	SIGCONT	Continue a stopped process.

By default, the bash shell ignores any `SIGQUIT` (3) and `SIGTERM` (15) signals that it receives (this is so that an interactive shell cannot be accidentally terminated). However, the bash shell does process any `SIGHUP` (1) and `SIGINT` (2) signals it receives.

If the bash shell receives a `SIGHUP` signal, it exits. Before it exits, however, it passes the `SIGHUP` signal to any processes started by the shell (such as your shell script). With a `SIGINT` signal, the shell is just interrupted. The Linux kernel stops giving the shell processing time on the CPU. When this happens, the shell passes the `SIGINT` signal to any processes started by the shell to notify them of the situation.

The shell passes these signals on to your shell script program for processing. The default behavior of shell scripts, however, is to ignore the signals, which may have an adverse effect on the operation of your script. To avoid this situation, you can program your script to recognize signals and perform commands to prepare the script for the consequences of the signal.

Generating Signals

The bash shell allows you to generate two basic Linux signals using key combinations on the keyboard. This feature comes in handy if you need to stop or pause a

runaway program.

Interrupting a Process

The Ctrl+C key combination generates a `SIGINT` signal, and sends it to any processes currently running in the shell. You can test this by running a command that normally takes a long time to finish, and pressing the Ctrl+C key combination:

```
$ sleep 100  
^C  
$
```

The Ctrl+C key combination simply just stops the current process running in the shell. The `sleep` command pauses the operation for the specified number of seconds. Normally, the command prompt wouldn't return until the timer has expired. By pressing the Ctrl+C key combination before the timer expires, you can cause the `sleep` command to terminate prematurely.

Pausing a Process

Instead of terminating a process, you can pause it in the middle of whatever it's doing. Sometimes this can be a dangerous thing (for example, if a script has a file lock open on a crucial system file), but often it allows you to peek inside what a script is doing without actually terminating the process.

The Ctrl+Z key combination generates a `SIGTSTP` signal, stopping any processes running in the shell. Stopping a

process is different than terminating the process, as stopping the process leaves the program still in memory and able to continue running from where it left off. In the “Job Control” section later on, you'll learn how to restart a process that's been stopped.

When you use the Ctrl+Z key combination, the shell informs you that the process has been stopped:

```
$ sleep 100
^Z
[1]+  Stopped                  sleep 100
$
```

The number in the square brackets is the *job number* assigned by the shell. The shell refers to each process running in the shell as a *job* and assigns each job a unique job number. It assigns the first process started job number 1, the second job number 2, and so on.

If you have a stopped job assigned to your shell session, bash will warn you if you try to exit the shell:

```
$ exit
logout
There are stopped jobs.
$
```

You can view the stopped job by using the `ps` command:

```
$ ps au
USER PID %CPU %MEM    VSZ   RSS TTY      STAT START   TIME COMMAND
rich 20560  0.0  1.2   2688 1624 pts/0    S   05:15   0:00 -bash
rich 20605  0.2  0.4   1564  552 pts/0    T   05:22   0:00 sleep 100
rich 20606  0.0  0.5   2584  740 pts/0    R   05:22   0:00 ps au
$
```

In the `STAT` column, the `ps` command shows the status of the stopped job as T. This indicates the command is either being traced or is stopped.

If you really want to exit the shell with the stopped job still active, just type the `exit` command again. The shell will exit, terminating the stopped job. Alternately, now that you know the PID of the stopped job, you can use the `kill` command to send a `SIGKILL` signal to terminate it:

```
$ kill -9 20605
$
[1]+  Killed                  sleep 100
$
```

When you kill the job, initially you won't get any response. However, the next time you do something that produces a shell prompt, you'll see a message indicating that the job was killed. Each time the shell produces a prompt, it also displays the status of any jobs that have changed states in the shell. After you kill a job, the next time you force the shell to produce a prompt it will display a message showing that the job was killed while running.

Trapping Signals

Instead of allowing your script to ignore signals, you can trap them when they appear and perform other commands. The `trap` command allows you to specify which Linux signals your shell script can watch for and intercept from the shell. If the script receives a signal listed in the `trap` command, it prevents it from being processed by the shell, and instead handles it locally.

The format of the `trap` command is:

```
trap commands signals
```

That's simple enough. On the `trap` command line, you just

list the commands you want the shell to execute, along with a space-separated list of signals you want to trap. You can specify the signals either by their numeric value or by their Linux signal name.

Here's a simple example of using the `trap` command to ignore `SIGINT` and `SIGTERM` signals:

```
$  
$ cat test1  
#!/bin/bash  
# testing signal trapping  
#  
trap "echo ` Sorry! I have trapped Ctrl-C`" SIGINT SIGTERM  
echo This is a test program  
count=1  
while [ $count -le 10 ]  
do  
    echo "Loop #$count"  
    sleep 5  
    count=$(( $count + 1 )  
done  
echo This is the end of the test program  
$  
$
```

The `trap` command used in this example displays a simple text message each time it detects either the `SIGINT` or `SIGTERM` signal. Trapping these signals makes this script impervious to the user attempting to stop the program by using the bash shell keyboard `Ctrl+C` command:

```
$  
$ ./test1  
This is a test program  
Loop #1  
Loop #2  
Loop #3  
^C Sorry! I have trapped Ctrl-C  
Loop #4  
Loop #5  
Loop #6  
Loop #7  
^C Sorry! I have trapped Ctrl-C  
Loop #8
```

```
Loop #9
Loop #10
This is the end of the test program
$
```

Each time the Ctrl+C key combination was used, the script executed the echo statement specified in the trap command instead of ignoring the signal and allowing the shell to stop the script.

Trapping a Script Exit

Besides trapping signals in your shell script, you can trap them when the shell script exits. This is a convenient way to perform commands just as the shell finishes its job.

To trap the shell script exiting, just add the `EXIT` signal to the `trap` command:

```
$ cat test2
#!/bin/bash
# trapping the script exit

trap "echo byebye" EXIT

count=1
while [ $count -le 5 ]
do
    echo "Loop #$count"
    sleep 3
    count=$(( $count + 1 ))
done
$
$ ./test2
Loop #1
Loop #2
Loop #3
Loop #4
Loop #5
byebye
$
```

When the script gets to the normal exit point, the trap is

triggered, and the shell executes the command you specify on the `trap` command line. The `EXIT` trap also works if you prematurely exit the script:

```
$ ./test2
Loop #1
Loop #2
^Cbyebye

$
```

When the `Ctrl+C` key combination is used to send a `SIGINT` signal, the script exits (because that signal isn't listed in the trap list), but before the script exits, the shell executes the `trap` command.

Removing a Trap

You can remove a set trap by using a dash as the command and a list of the signals you want to return to normal behavior:

```
$ cat test3
#!/bin/bash
# removing a set trap

trap "echo byebye" EXIT

count=1
while [ $count -le 5 ]
do
    echo "Loop #$count"
    sleep 3
    count=$(( $count + 1 ))
done
trap - EXIT
echo "I just removed the trap"
$
$ ./test3
Loop #1
Loop #2
Loop #3
```

```
Loop #4
Loop #5
I just removed the trap
$
```

Once the signal trap is removed, the script ignores the signals. However, if a signal is received before the trap is removed, the script processes it per the `trap` command:

```
$ ./test3
Loop #1
Loop #2
^Cbyebye

$
```

In this example, a Ctrl+C key combination was used to terminate the script prematurely. Because the script was terminated before the trap was removed, the script executed the command specified in the trap.

Running Scripts in Background Mode

There are times when running a shell script directly from the command line interface is inconvenient. Some scripts can take a long time to process, and you may not want to tie up the command line interface waiting. While the script is running, you can't do anything else in your terminal session. Fortunately, there's a simple solution to that problem.

When you use the `ps` command, you see a whole bunch of different processes running on the Linux system. Obviously, all of these processes are not running on your terminal

monitor. This is called running processes in the *background*. In background mode, a process runs without being associated with a `STDIN`, `STDOUT`, and `STDERR` on a terminal session (see Chapter 14).

You can exploit this feature with your shell scripts as well, allowing them to run behind the scenes and not lock up your terminal session. The following sections describe how to run your scripts in background mode on your Linux system.

Running in the Background

Running a shell script in background mode is a fairly easy thing to do. To run a shell script in background mode from the command line interface, just place an ampersand symbol after the command:

```
$ ./test1 &
[1] 1976
$ This is a test program
Loop #1
Loop #2
ls /home/user/Desktop
gnome-screenshot.desktop  konsole.desktop      virtualbox.desktop
gnome-terminal.desktop   ksnapshot.desktop
$ Loop #3
Loop #4
...
$
```

When you place the ampersand symbol after a command, it separates the command from the bash shell and runs it as a separate background process on the system. The first thing that displays is the line:

```
[1] 1976
```

The number in the square brackets is the job number assigned to the background process by the shell. The next number is the Process ID (PID) the Linux system assigns to the process. Every process running on the Linux system must have a unique PID.

As soon as the system displays these items, a new command line interface prompt appears. You are returned back to the shell, and the command you executed runs safely in background mode.

At this point, you can enter new commands at the prompt (as shown in the example). However, while the background process is still running, it still uses your terminal monitor for `STDOUT` and `STDERR` messages. You'll notice from the example that the output from the `test1` script appears in the output intermixed with any other commands that are run from the shell.

When the background process finishes, it displays a message on the terminal:

```
[1]+  Done                  ./test1
```

This shows the job number and the status of the job (Done), along with the command used to start the job.

Running Multiple Background Jobs

You can start any number of background jobs at the same time from the command line prompt:

```
$  
$ ./test1 &  
[3] 2174  
$ This is Test Program 1
```



```

$ ./test2 &
[4] 2176
$ I am Test Program 2

$ ./test3 &
[5] 2178
$ Well this is Test Program

$ ./test4 &
[6] 2180
$ This is Test Program 4

$

```

Each time you start a new job, the Linux system assigns it a new job number and PID. You can see that all of the scripts are running using the `ps` command:

```

$ ps au

```

USER	PID	%CPU	%MEM	VSZ	RSS	TTY	STAT	START	TIME	COMMAND
...										
user	1826	0.0	0.3	6704	3408	pts/0	Ss	14:07	0:00	bash
user	2174	0.0	0.1	4860	1076	pts/0	S	15:23	0:00	/bin/bash ./test1
user	2175	0.0	0.0	3884	504	pts/0	S	15:23	0:00	sleep 300
user	2176	0.0	0.1	4860	1068	pts/0	S	15:23	0:00	/bin/bash ./test2
user	2177	0.0	0.0	3884	508	pts/0	S	15:23	0:00	sleep 300
user	2178	0.0	0.1	4860	1068	pts/0	S	15:23	0:00	/bin/bash ./test3
user	2179	0.0	0.0	3884	504	pts/0	S	15:23	0:00	sleep 300
user	2180	0.0	0.1	4860	1068	pts/0	S	15:23	0:00	/bin/bash ./test4
user	2181	0.0	0.0	3884	504	pts/0	S	15:23	0:00	sleep 300
user	2182	0.0	0.1	4592	1100	pts/0	R+	15:24	0:00	ps au

```

$

```

Each of the background processes you start appears in the `ps` command output listing of running processes. If all of the processes display output in your terminal session, things can get pretty messy pretty quickly. Fortunately, there's a simple way to solve that problem, which we discuss in the next section.

Exiting the Terminal

You must be careful when using background processes from a terminal session. Notice in the output from the `ps` command that each of the background processes is tied to the terminal session (pts/0) terminal. If the terminal session exits, the background process also exits.

Some terminal emulators warn you if you have any running background processes associated with the terminal, while others don't. If you want your script to continue running in background mode after you have logged off the console, there's something else you need to do. The next section discusses that process.

Running Scripts Without a Console

There will be times when you want to start a shell script from a terminal session and then let the script run in background mode until it finishes, even if you exit the terminal session. You can do this by using the `nohup` command.

The `nohup` command runs another command blocking any `SIGHUP` signals that are sent to the process. This prevents the process from exiting when you exit your terminal session.

The format used for the `nohup` command is as follows:

```
$ nohup ./test1 &
[1] 19831
$ nohup: ignoring input and appending output to 'nohup.out'
$
```

As with a normal background process, the shell assigns the command a job number, and the Linux system assigns a PID number. The difference is that when you use the `nohup` command, the script ignores any `SIGHUP` signals sent by the terminal session if you close the session.

Because the `nohup` command disassociates the process from the terminal, the process loses the `STDOUT` and `STDERR` output links. To accommodate any output generated by the command, the `nohup` command automatically redirects `STDOUT` and `STDERR` messages to a file, called `nohup.out`.

The `nohup.out` file contains all of the output that would normally be sent to the terminal monitor. After the process finishes running, you can view the `nohup.out` file for the output results:

```
$ cat nohup.out
This is a test program
Loop #1
Loop #2
Loop #3
Loop #4
Loop #5
Loop #6
Loop #7
Loop #8
Loop #9
Loop #10
This is the end of the test program
$
```

The output appears in the `nohup.out` file just as if the process ran on the command line!

Caution

If you run another command using `nohup`, the output is

appended to the existing `nohup.out` file. Be careful when running multiple commands from the same directory, as all of the output will be sent to the same `nohup.out` file, which can get confusing.

Job Control

Earlier in this chapter, you saw how to use the key combination to stop a job running in the shell. After you stop a job, the Linux system lets you either kill or restart it. You can kill the process by using the `kill` command. Restarting a stopped process requires that you send it a `SIGCONT` signal.

The function of starting, stopping, killing, and resuming jobs is called *job control*. With job control, you have full control over how processes run in your shell environment.

This section describes the commands to use to view and control jobs running in your shell.

Viewing Jobs

The key command for job control is the `jobs` command. The `jobs` command allows you to view the current jobs being handled by the shell:

```
$ cat test4
#!/bin/bash
# testing job control

echo "This is a test program $$"
count=1
while [ $count -le 10 ]
do
```

```

echo "Loop #\$count"
sleep 10
count=$((count + 1))
done
echo "This is the end of the test program"
$
$ ./test4
This is a test program 29011
Loop #1
^Z

[1]+  Stopped                  ./test4
$
$ ./test4 > test4out &
[2] 28861
$
$ jobs
[1]+  Stopped                  ./test4
[2]-  Running                  ./test4 >test4out &
$

```

The script uses the `$$` variable to display the PID that the Linux system assigns to the script; then it goes into a loop, sleeping for 10 seconds at a time for each iteration. In the example, the first script is started from the command line interface and then stopped using the Ctrl+Z key combination. Next, another job is started as a background process, using the ampersand symbol. To make life a little easier, the output of that script is redirected to a file so that it doesn't appear on the screen.

After the two jobs were started, we used the `jobs` command to view the jobs assigned to the shell. The `jobs` command shows both the stopped and the running jobs, along with their job numbers and the commands used in the jobs.

The `jobs` command uses a few different command line parameters, as shown in [Table 15.2](#).

[Table 15.2](#) The jobs Command Parameters

Parameter	Description
-l	List the PID of the process along with the job number.
-n	List only jobs that have changed their status since the last notification from the shell.
-p	List only the PIDs of the jobs.
-r	List only the running jobs.
-s	List only stopped jobs.

You probably noticed the plus and minus signs in the `jobs` command output. The job with the plus sign is considered the default job. It would be the job referenced by any job control commands if a job number wasn't specified in the command line. The job with the minus sign is the job that would become the default job when the current default job finishes processing. There will only be one job with the plus sign and one job with the minus sign at any time, no matter how many jobs are running in the shell.

The following is an example showing how the next job in line takes over the default status, when the default job is removed:

```
$ ./test4
This is a test program 29075
Loop #1
^Z
[1]+  Stopped                  ./test4
$
$ ./test4
This is a test program 29090
Loop #1
^Z
[2]+  Stopped                  ./test4
$
$ ./test4
This is a test program 29105
Loop #1
```

```

2
[3]+  Stopped                  ./test4
$
$ jobs -l
[1] 29075 Stopped              ./test4
[2]- 29090 Stopped              ./test4
[3]+ 29105 Stopped              ./test4
$
$ kill -9 29105
$
$ jobs -l
[1]- 29075 Stopped              ./test4
[2]+ 29090 Stopped              ./test4
$

```

In this example, three separate processes were started and then stopped. The `jobs` command listing shows the three processes and their status. Note that the default process (the one listed with the plus sign) is the last process started.

Then the `kill` command was issued to send a `SIGKILL` signal to the default process. In the next `jobs` listing, the job that previously had the minus sign is now the default job.

Restarting Stopped Jobs

Under bash job control, you can restart any stopped job as either a background process or a foreground process. A foreground process takes over control of the terminal you're working on, so be careful about using that feature.

To restart a job in background mode, use the `bg` command, along with the job number:

```

$ bg 2
[2]+ ./test4 &
Loop #2
$ Loop #3
Loop #4

```

```

$ jobs
[1]+  Stopped                  ./test4
[2]-  Running                  ./test4 &
$ Loop #6
Loop #7
Loop #8
Loop #9
Loop #10
This is the end of the test program
[2]-  Done                     ./test4
$

```

Because the job was restarted in background mode, the command line interface prompt appears, allowing other commands to be entered. The output from the `jobs` command now shows that the job is indeed running (as you can tell from the output now appearing on the monitor).

To restart a job in foreground mode, use the `fg` command, along with the job number:

```

$ jobs
[1]+  Stopped                  ./test4
$ fg 1
./test4
Loop #2
Loop #3

```

Since the job is running in foreground mode, the command line interface prompt does not appear until the job finishes.

Being Nice

In a multitasking operating system (which Linux is), the kernel is responsible for assigning CPU time for each process running on the system. Only one process at a time

can actually be running in a CPU, so the kernel assigns CPU time to each process in turn.

By default, all processes started from the shell have the same *scheduling priority* on the Linux system. The scheduling priority is the amount of CPU time the kernel assigns to the process relative to the other processes.

The scheduling priority is an integer value, from -20 (the highest priority) to + 20 (the lowest priority). By default, the bash shell starts all processes with a priority of 0.

Tip

It's confusing to remember that -20, the lowest value, is the highest priority and 19, the highest value, is the lowest priority. Just remember the phrase, "Nice guys finish last." The "nicer" or higher you are in value, the lower your chance of getting the CPU.

This means that a simple script that requires only a little bit of processing time gets the same CPU time slices as a complex mathematical algorithm that can take hours to run.

Sometimes you want to change the priority of a specific command, either lowering its priority so that it doesn't take as much processing power from the CPU or giving it a higher priority so that it gets more processing time. You can do this by using the `nice` command.

The nice Command

The `nice` command allows you to set the scheduling priority of a command as you start it. To make a command run with less priority, just use the `-n` command line option for `nice` to specify a new priority level:

```
$ nice -n 10 ./test4 > test4out &
[1] 29476
$ ps al
 F      UID      PID      PPID     PRI      NI  WCHAN      STAT  TTY      TIME COMMAND
100    501    29459    29458     12       0  wait4      S      pts/0    0:00 -bash
000    501    29476    29459     15      10  wait4      SN     pts/0    0:00 /bin/bash
000    501    29490    29476     15      10  nanosl     SN     pts/0    0:00 sleep 10
000    501    29491    29459     14       0  -          R      pts/0    0:00 ps al
$
```

The `nice` command causes the script to run at a lower priority. However, if you try to increase the priority of one of your commands, you might be in for a surprise:

```
$ nice -n -10 ./test4 > test4out &
[1] 29501
$ nice: cannot set priority: Permission denied

[1]+  Exit 1                  nice -n -10 ./test4 >test4out
$
```

The `nice` command prevents normal system users from increasing the priority of their commands. This is a safety feature to prevent a user from starting all of his or her commands as high priority.

The renice Command

Sometimes you'd like to change the priority of a command that's already running on the system. That's what the `renice` command is for. It allows you to specify the PID of a running process to change its priority:

```
$ ./test4 > test4out &
[1] 29504
```

```

$ ps al
  F   UID     PID   PPID  PRI   NI  WCHAN   STAT  TTY          TIME COMMAND
100   501   29459   29458   12    0  wait4    S     pts/0        0:00 -bash
000   501   29504   29459    9    0  wait4    S     pts/0        0:00 /bin/bash .
000   501   29518   29504    9    0  nanosl   S     pts/0        0:00 sleep 10
000   501   29519   29459   14    0  -        R     pts/0        0:00 ps al

$ renice 10 -p 29504
29504: old priority 0, new priority 10

$ ps al
  F   UID     PID   PPID  PRI   NI  WCHAN   STAT  TTY          TIME COMMAND
100   501   29459   29458   16    0  wait4    S     pts/0        0:00 -bash
000   501   29504   29459   14   10  wait4    SN    pts/0        0:00 /bin/bash .
000   501   29535   29504    9    0  nanosl   S     pts/0        0:00 sleep 10
000   501   29537   29459   14    0  -        R     pts/0        0:00 ps al

$

```

The `renice` command automatically updates the scheduling priority of the running process. As with the `nice` command, the `renice` command has some limitations:

- You can only `renice` processes that you own.
- You can only `renice` your processes to a lower priority.
- The root user can `renice` any process to any priority.

If you want to fully control running processes, you must be logged in as the root account or use the `sudo` command

Running Like Clockwork

I'm sure that, as you start working with scripts, there will be a situation in which you will want to run a script at a preset time, usually at a time when you're not there. The Linux system provides a couple of ways to run a script at a preselected time: the `at` command and the `cron` table. Each method uses a different technique for scheduling when and how often to run scripts. The following sections describe

each of these methods.

Scheduling a Job Using the `at` Command

The `at` command allows you to specify a time when the Linux system will run a script. The `at` command submits a job to a queue with directions on when the shell should run the job. The `at` daemon, `atd`, runs in the background and checks the job queue for jobs to run. Most Linux distributions start this daemon automatically at boot time.

The `atd` daemon checks a special directory on the system (usually `/var/spool/at`) for jobs submitted using the `at` command. By default, the `atd` daemon checks this directory every 60 seconds. When a job is present, the `atd` daemon checks the time the job is set to be run. If the time matches the current time, the `atd` daemon runs the job.

The following sections describe how to use the `at` command to submit jobs to run and how to manage jobs.

The `at` Command Format

The basic `at` command format is pretty simple:

```
at [-f filename] time
```

By default, the `at` command submits input from `STDIN` to the queue. You can specify a file name used to read commands (your script file) using the `-f` parameter.

The time parameter specifies when you want the Linux

system to run the job. You can get pretty creative with how you specify the time. The `at` command recognizes lots of different time formats:

- A standard hour and minute, such as 10:15
- An ~A.M./~P.M. indicator, such as 10:15~P.M.
- A specific named time, such as now, noon, midnight, or teatime (4~P.M.)

If you specify a time that has already past, the `at` command runs the job at that time on the next day.

In addition to specifying the time to run the job, you can also include a specific date, using a few different date formats:

- A standard date format, such as MMDDYY, MM/DD/YY, or DD.MM.YY
- A text date, such as Jul 4 or Dec 25, with or without the year
- You can also specify a time increment:
 - Now + 25 minutes
 - 10:15~P.M. tomorrow
 - 10:15 + 7 days

When you use the `at` command, the job is submitted into a *job queue*. The job queue holds the jobs submitted by the `at` command for processing. There are 26 different job queues available for different priority levels. Job queues are referenced using lowercase letters, a through z.

Note

A few years ago, the `batch` command was another method that allowed a script to be run at a later time. The `batch` command was unique because you could schedule a script to run when the system was at a lower usage level. However, nowadays, the `batch` command is just simply a script, `/usr/bin/batch`, that calls the `at` command and submits your job to the `batch` queue.

The higher alphabetically the job queue, the lower the priority (higher `nice` value) the job will run under. By default, `at` jobs are submitted to the `at` job queue. If you want to run a job at a higher priority, you can specify a different queue letter using the `-q` parameter.

Retrieving Job Output

When the job runs on the Linux system, there's no monitor associated with the job. Instead, the Linux system uses the e-mail address of the user who submitted the job as `STDOUT` and `STDERR`. Any output destined to `STDOUT` or `STDERR` is mailed to the user via the mail system.

Here's a simple example of using the `at` command to schedule a job to run:

```
$
$ cat test5
#!/bin/bash
#
# testing the at command
#
echo This script ran at `date`
echo This is the end of the script >&2
$
$ date
```

```
Mon Oct 18 14:38:17 DT 2010
$
$ at -f test5 14:39
warning: commands will be executed using /bin/sh
job 57 at Mon Oct 18 14:39:00 2010
$
$ mail
"/var/mail/user": 1 message 1 new
>N      1 user              Mon Oct 18 14:39   15/538   Output from your job
& 1
Date: Mon, 18 Oct 2010 14:39:00 -0400
Subject: Output from your job          57
To: user@user-desktop
From: user <user@user-desktop>

This script ran at Mon Oct 18 14:39:00 EDT 2010
This is the end of the script

& exit
$
```

The `at` command produces a warning message, indicating what shell the system uses to run the script, `/bin/sh`, along with the job number assigned to the job and the time the job is scheduled to run.

Caution

In most Linux distributions, the default shell assigned to `/bin/sh` is the bash shell. However, Ubuntu uses the dash shell for its default shell. See Chapter 22 for more information on the dash shell.

When the job completes, nothing appears on the monitor, but the system generates an e-mail message. The e-mail message shows the output generated by the script. If the script doesn't produce any output, it won't generate an e-mail message, by default. You can change that by using the

-m option in the `at` command. This generates an e-mail message, indicating the job completed, even if the script doesn't generate any output.

Listing Pending Jobs

The `atq` command allows you to view what jobs are pending on the system:

```
$
$ at -f test5 15:05
warning: commands will be executed using /bin/sh
job 58 at Mon Oct 18 15:05:00 2010
$
$ at -f test5 15:10
warning: commands will be executed using /bin/sh
job 59 at Mon Oct 18 15:10:00 2010
$
$ at -f test5 15:15
warning: commands will be executed using /bin/sh
job 60 at Mon Oct 18 15:15:00 2010
$
$ at -f test5 15:20
warning: commands will be executed using /bin/sh
job 61 at Mon Oct 18 15:20:00 2010
$$ atq
61      Mon Oct 18 15:20:00 2010 a user
58      Mon Oct 18 15:05:00 2010 a user
59      Mon Oct 18 15:10:00 2010 a user
60      Mon Oct 18 15:15:00 2010 a user
$
```

The job listing shows the job number, the date and time the system will run the job, and the job queue the job is stored in.

Removing Jobs

Once you know the information about what jobs are pending in the job queues, you can use the `atrm` command

to remove a pending job:

```
$  
$ atq  
59      Mon Oct 18 15:10:00 2010 a user  
60      Mon Oct 18 15:15:00 2010 a user  
$  
$ atrm 59  
$  
$ atq  
60      Mon Oct 18 15:15:00 2010 a user  
$
```

Just specify the job number you want to remove. You can only remove jobs that you submit for execution. You can't remove jobs submitted by others.

Scheduling Regular Scripts

Using the `at` command to schedule a script to run at a preset time is great, but what if you need that script to run at the same time every day or once a week or once a month? Instead of having to continually submit `at` jobs, you can use another feature of the Linux system.

The Linux system uses the `cron` program to allow you to schedule jobs that need to run on a regular basis. The `cron` program runs in the background and checks special tables, called *cron tables*, for jobs that are scheduled to run.

The cron Table

The `cron` table uses a special format for allowing you to specify when a job should be run. The format for the `cron` table is:

```
min hour dayofmonth month dayofweek command
```

The `cron` table allows you to specify entries as specific values, ranges of values (such as 1–5) or as a wildcard character (the asterisk). For example, if you want to run a command at 10:15 on every day, you would use the `cron` table entry of:

```
15 10 * * * command
```

The wildcard character used in the *dayofmonth*, *month*, and *dayofweek* fields indicates that `cron` will execute the command every day of every month at 10:15. To specify a command to run at 4:15 PM every Monday, you would use the following:

```
15 16 * * 1 command
```

You can specify the *dayofweek* entry as either a three-character text value (mon, tue, wed, thu, fri, sat, sun) or as a numeric value, with 0 being Sunday and 6 being Saturday.

Here's another example: to execute a command at 12 noon on the first day of every month, you would use the following format:

```
00 12 1 * * command
```

The *dayofmonth* entry specifies a date value (1–31) for the month.

Note

The astute reader might be wondering just how you would be able to set a command to execute on the last day of every month because you can't set the *dayofmonth* value to cover every month. This problem has plagued Linux and Unix programmers, and has

spawned quite a few different solutions. A common method is to add an `if-then` statement that uses the `date` command to check if tomorrow's date is 01:

```
00 12 * * * if [ `date +%d -d tomorrow` =  
01 ] ; then ; command
```

This will check every day at 12 noon to see if it's the last day of the month, and if so, `cron` will run the command.

The command list must specify the full pathname of the command or shell script to run. You can add any command line parameters or redirection symbols you like, as a regular command line:

```
15 10 * * * /home/rich/test4 > test4out
```

The `cron` program runs the script using the user account that submitted the job. Thus, you must have the proper permissions to access the command and output files specified in the command listing.

Building the cron Table

Each system user can have their own `cron` table (including the root user) for running scheduled jobs. Linux provides the `crontab` command for handling the `cron` table. To list an existing `cron` table, use the `-l` parameter:

```
$ crontab -l  
no crontab for rich  
$
```

By default, each user's `cron` table file doesn't exist. To add entries to your `cron` table, use the `-e` parameter. When you do that, the `crontab` command starts a text editor (see

Chapter 9) with the existing `cron` table (or an empty file if it doesn't yet exist).

cron Directories

When you create a script that has less precise execution time needs, it is easier to use one of the pre-configured `cron` script directories. There are four basic directories: hourly, daily, monthly, and weekly.

```
$  
$ ls /etc/cron.*ly  
/etc/cron.daily:  
0anacron          aptitude          exim4-base        man-db            sysstat  
apt               bsdmainutils     logrotate         popularity-contest  
appport           dpkg              mlocate          standard  
/etc/cron.hourly:  
  
/etc/cron.monthly:  
0anacron  
  
/etc/cron.weekly:  
0anacron apt-xapian-index man-db  
$
```

Thus, if you have a script that needs to be run one time per day, just copy the script to the daily directory and `cron` executes it each day.

The anacron Program

The only problem with the `cron` program is that it assumes that your Linux system is operational 24 hours a day, 7 days a week. Unless you're running Linux in a server environment, this may not necessarily be true.

If the Linux system is turned off at the time a job is scheduled to run in the `crontable`, the job won't run. The `cron`

program doesn't retroactively run missed jobs when the system is turned back on. To resolve this issue, many Linux distributions also include the `anacron` program.

If `anacron` determines that a job has missed a scheduled running, it runs the job as soon as possible. This means that if your Linux system is turned off for a few days, when it starts back up any jobs scheduled to run during the time it was off are automatically run.

This is a feature that's often used for scripts that perform routine log maintenance. If the system is always off when the script should run, the log files would never get trimmed and could grow to undesirable sizes. With `anacron`, you're guaranteed that the log files will be trimmed at least each time the system is started.

The `anacron` program deals only with programs located in the `cron` directories, such as `/etc/cron.monthly`. It uses timestamps to determine if the jobs have been run at the proper scheduled interval. A timestamp file exists for each `cron` directory and is located in `/var/spool/anacron`:

```
$
$ sudo cat /var/spool/anacron.monthly
20101123
$
```

The `anacron` program has its own table (usually located at `/etc/anacrontab`) to check the job directories:

```
$
$ cat /etc/anacrontab
# /etc/anacrontab: configuration file for anacron

# See anacron(8) and anacrontab(5) for details.

SHELL=/bin/sh
```

```
PATH=/usr/local/sbin:/usr/local/bin:/sbin:/bin:/usr/sbin:/usr/bin
# These replace cron's entries
1 5 cron.daily nice run-parts --report /etc/cron.daily
7 10 cron.weekly nice run-parts --report /etc/cron.weekly
@monthly 15 cron.monthly nice run-parts --report /etc/cron.monthly
$
```

The basic format of the `anacron` table is slightly different from that of the `cron` table:

period delay identifier command

The period entry defines how often the jobs should be run, specified in days. The `anacron` program uses this entry to check against the jobs' timestamp file. The delay entry specifies how many minutes after the system starts the `anacron` program should run missed scripts. The command entry contains the `run-parts` program and a `cron` script directory name. The `run-parts` program is responsible for running any script in the directory passed to it.

Notice that `anacron` does not run the scripts located in `/etc/cron.hourly`. This is because the `anacron` program does not deal with scripts that have execution time needs of less than daily.

The identifier entry is a unique non-blank character string—for example, `cron-weekly`. It is used to uniquely identify the job in log messages and error e-mails.

Start at the Beginning

The last method of starting shell scripts is to have your script run automatically either as soon as the Linux system boots or whenever a user starts a new bash shell session.

Starting scripts at boot time is usually reserved for special scripts that perform system functions, such as configuring a network interface or starting a server process. However, if you're a Linux system administrator, it's possible that you'll need to perform a function every time the Linux system boots, such as resetting a custom log file or starting a custom application.

The ability to run a script every time a user starts a new bash shell (even just when a specific user starts a bash shell) also can come in handy. There are times when you want to set shell features for a shell session or just ensure that a specific file has been set.

This section describes how to configure your Linux system to run your scripts either at boot time or each time a new bash shell starts.

Starting Your Scripts at Boot Time

Before you can get your shell script to start at boot time, you'll need to know a little bit about how the Linux boot process works. There's a specific order that Linux uses to start scripts at boot time, and knowing that process can help you in getting your script to perform the way you want it.

The Boot Process

After you turn on your Linux system, the Linux kernel loads into memory and runs. The first thing it does is start the

UNIX System V init process or the Upstart init process, depending upon your distribution and its version. This process is then responsible for starting all other processes on the Linux system.

System V init Process

As part of the boot process, the System V init process reads the `/etc/inittab` file. The `inittab` file lists the system's *run level*. Different Linux run levels start different programs and scripts. [Table 15.3](#) lists the Linux run levels for Red Hat–based distributions.

Table 15.3 The Linux Run Levels—Red Hat–Based Distributions

Run Level	Description
0	Halt
1	Single-user mode
2	Multi-user mode, usually without networking support
3	Full multi-user mode, with networking
4	User definable
5	Multi-user mode, with networking and a graphical XWindow session
6	Reboot

It is common for Linux distributions to use the run levels listed in [Table 15.3](#) to grant finer control over what services are started on the system. However, Debian-based distributions, such as Ubuntu and Linux Mint, do not

distinguish a difference between run levels 2 through 5, as shown in [Table 15.4](#)

Table 15.4 The Linux Run Levels—Debian-Based Distributions

Run Level	Description
0	Halt
1	Single-user mode
2 - 5	Multi-user mode, with networking and a graphical XWindow session
6	Reboot

The Ubuntu distribution does not even have a `/etc/inittab` file to set run levels. By default, Ubuntu goes to run level 2. You must create a `/etc/inittab` file on an Ubuntu distribution, if you want to change the default run level.

Each run level defines what scripts the System V init process starts or stops. These *startup script* are shell scripts that start applications, providing the necessary environment variables for them to run.

This is the part of the System V init based boot process where things start to get a little fuzzy, mainly because different Linux distributions place startup scripts in slightly different locations. Some distributions place startup scripts in the `/etc/rc#.d` directory, where `#` stands for the run level. Others use the `/etc/init.d` directory, and still others use the `/etc/init.d/rc.d` directory. Usually a quick glance in your `/etc` directory structure can easily determine what format your

distribution uses.

Upstart init Process

The Upstart init process is a newer standard for managing service processes that many Linux distributions are now moving toward. Upstart is not focused on system run levels, but is instead focused on events, such as a system boot. In Upstart, a system boot is called a *startup event*.

Upstart uses files located in either the `/etc/event.d` or the `/etc/init` directory to start processes, depending upon your distribution and its version. In order to maintain backward compatibility, many Upstart implementations still call the older System V init scripts from `/etc/init.d` and/or `/etc/rc#.d`.

The Upstart init process standard is still highly evolving at this point in time. For the scripts you create and want to start during a system startup event, it would be prudent to use the older System V init process method. This method is described in the next section.

Defining Your Scripts

It's best not to mess with the individual startup script files in your Linux distribution. Often distributions provide tools to automatically build these scripts as you add server applications, and manually changing these scripts can cause problems.

Instead, most Linux distributions provide a local startup file specifically to allow the system administrator to enter

scripts to run at boot time. Of course, the name and location of this file are different in different Linux distributions. [Table 15.5](#) identifies the location of the startup file in five popular Linux distributions.

Table 15.5 The Linux Local Startup File Locations

Distribution	File location
debian	/etc/init.d/rc.local
Fedora	/etc/rc.d/rc.local
Mandriva	/etc/rc.local
openSuse	/etc/init.d/boot.local
Ubuntu	/etc/rc.local

Inside the local startup file, you can either specify specific commands and statements, or enter any scripts you want started at boot time. Remember that if you use a script, you'll need to specify the full pathname for the script so that the system can find it at boot time.

Caution

Different Linux distributions also execute the local startup script at different points in the boot process. Sometimes the script is run before things such as network support have been started. Consult your specific Linux distribution documentation to determine when the local startup script is run in your distribution.

Starting with a New Shell

Each user's home directory contains two files that the bash shell uses to automatically start scripts and set environment variables:

- The `.bash_profile` file
- The `.bashrc` file

The bash shell runs the `.bash_profile` file when a new shell is run as a result of a new login. Place any scripts that you want run at login time in this file.

The bash shell runs the `.bashrc` file any time a new shell is started, including when a new login occurs. You can test this by adding a simple echo statement to the `.bashrc` file in your home directory and then starting a new shell:

```
$ bash
This is a new shell!!
$
```

If you want to run a script for every user on the system, most Linux distributions provide the `/etc/bashrc` file (note that there's no period in front of the `bashrc` filename). The bash shell executes the statements in this file every time any user on the system starts a new bash shell.

Summary

The Linux system allows you to control your shell scripts by using signals. The bash shell accepts signals and passes them on to any process running under the shell process. Linux signals allow you to easily kill a runaway process or temporarily pause a long-running process.

You can use the `trap` statement in your scripts to catch signals and perform commands. This feature provides a simple way to control whether a user can interrupt your script while it's running.

By default, when you run a script in a terminal session shell, the interactive shell is suspended until the script completes. You can cause a script or command to run in background mode by adding an ampersand sign (&) after the command name. When you run a script or command in background mode, the interactive shell returns, allowing you to continue entering more commands. Any background processes run using this method are still tied to the terminal session. If you exit the terminal session, the background processes also exit.

To prevent this from happening, use the `nohup` command. This command intercepts any signals intended for the command that would stop it—for example, when you exit the terminal session. This allows scripts to continue running in background mode even if you exit the terminal session.

When you move a process to background mode, you can still control what happens to it. The `jobs` command allows you to view processes started from the shell session. Once you know the job ID of a background process, you can use the `kill` command to send Linux signals to the process or use the `fg` command to bring the process back to the foreground in the shell session. You can suspend a running foreground process by using the `Ctrl+Z` key combination and then place it back in background mode, using the `bg`

command.

The `nice` and `renice` commands allow you to change the priority level of a process. By giving a process a lower priority, you allow the CPU to allocate less time to it. This comes in handy when running long processes that can take lots of CPU time.

In addition to controlling processes while they're running, you can also determine when a process starts on the system. Instead of running a script directly from the command line interface prompt, you can schedule the process to run at an alternative time. There are several different ways to accomplish this. The `at` command enables you to run a script once at a preset time. The `cron` program provides an interface that can run scripts at a regularly scheduled interval.

Finally, the Linux system provides script files for you to use for scheduling your scripts to run either at system boot time or whenever a user starts a new bash shell. Distributions have different files that allow you to list scripts that start each time the system boots. This permits system administrators to run special scripts for system maintenance at boot time. Similarly, the `.bash_profile` and `.bashrc` files are located in every user's home directory to provide a location to place scripts and commands that run with a new shell. The `.bash_profile` file runs scripts each time a user logs in to the system, and the `.bashrc` file runs scripts on each new shell instance.

In the next chapter, we'll look at how to write script

functions. Script functions allow you to write code blocks once and then use them in multiple locations throughout your script.

Part III

Advanced Shell Scripting

In This Part

Chapter 16: Creating Functions

Chapter 17: Writing Scripts for Graphical Desktops

Chapter 18: Introducing sed and gawk

Chapter 19: Regular Expressions

Chapter 20: Advanced sed

Chapter 21: Advanced gawk

Chapter 22: Working with Alternative Shells

Chapter 16

Creating Functions

In This Chapter

- Basic script functions
- Returning a value
- Using variables in functions
- Array and variable functions
- Function recursion
- Creating a library
- Using functions on the command line

Often while writing shell scripts you'll find yourself using the same code in multiple locations. If it's just a small code snippet, it's usually not that big of a deal. However, rewriting large chunks of code multiple times in your shell script can get tiring. The bash shell provides a way to help you out by supporting user-defined functions. You can encapsulate your shell script code into a function and then use it as many times as you want anywhere in your script. This chapter walks you through the process of creating your own shell script functions and demonstrates how to use them in other shell script applications.

Basic Script Functions

As you start writing more complex shell scripts, you'll find yourself reusing parts of code that perform specific tasks. Sometimes it's something simple, such as displaying a text message and retrieving an answer from the script users. Other times it's a complicated calculation that's used multiple times in your script as part of a larger process.

In each of these situations, it can get tiresome writing the same blocks of code over and over again in your script. It would be nice to just write the block of code once and then be able to refer to that block of code anywhere in your script without having to rewrite it.

The bash shell provides a feature allowing you to do just that. *Functions* are blocks of script code that you assign a name to and then reuse anywhere in your code. Anytime you need to use that block of code in your script, you simply use the function name you assigned it (referred to as *calling* the function). This section describes how to create and use functions in your shell scripts.

Creating a Function

There are two formats you can use to create functions in bash shell scripts. The first format uses the keyword `function`, along with the function name you assign to the block of code:

```
function name {  
    commands  
}
```

The *name* attribute defines a unique name assigned to the function. Each function you define in your script must be assigned a unique name.

The *commands* are one or more bash shell commands that make up your function. When you call the function, the bash shell executes each of the commands in the order they appear in the function, just as in a normal script.

The second format for defining a function in a bash shell script more closely follows how functions are defined in other programming languages:

```
name() {  
    commands  
}
```

The empty parentheses after the function name indicate that you're defining a function. The same naming rules apply in this format as in the original shell script function format.

Using Functions

To use a function in your script, specify the function name on a line, just as you would any other shell command:

```
$ cat test1  
#!/bin/bash  
# using a function in a script  
  
function func1 {  
    echo "This is an example of a function"  
}
```

```

count=1
while [ $count -le 5 ]
do
    func1
    count=$(( $count + 1 ))
done

echo "This is the end of the loop"
func1
echo "Now this is the end of the script"
$
$ ./test1
This is an example of a function
This is an example of a function
This is an example of a function
This is an example of a function
This is an example of a function
This is the end of the loop
This is an example of a function
Now this is the end of the script
$

```

Each time you reference the `func1` function name, the bash shell returns to the `func1` function definition and executes any commands you defined there.

The function definition doesn't have to be the first thing in your shell script, but be careful. If you attempt to use a function before it's defined, you'll get an error message:

```

$ cat test2
#!/bin/bash
# using a function located in the middle of a script

count=1
echo "This line comes before the function definition"

function func1 {
    echo "This is an example of a function"
}

while [ $count -le 5 ]
do
    func1
    count=$(( $count + 1 ))
done
echo "This is the end of the loop"

```

```

func2
echo "Now this is the end of the script"

function func2 {
    echo "This is an example of a function"
}
$
$ ./test2
This line comes before the function definition
This is an example of a function
This is an example of a function
This is an example of a function
This is an example of a function
This is an example of a function
This is the end of the loop
./test2: func2: command not found
Now this is the end of the script
$

```

The first function, `func1`, was defined after a couple of statements in the script, which is perfectly fine. When the `func1` function was used in the script, the shell knew where to find it.

However, the script attempted to use the `func2` function before it was defined. Because the `func2` function wasn't defined, when the script reached the place where we used it, it produced an error message.

You also need to be careful about your function names. Remember, each function name must be unique, or you'll have a problem. If you redefine a function, the new definition will override the original function definition, without producing any error messages:

```

$ cat test3
#!/bin/bash
# testing using a duplicate function name

function func1 {
echo "This is the first definition of the function name"
}

```

```
func1

function func1 {
    echo "This is a repeat of the same function name"
}

func1

echo "This is the end of the script"
$
$ ./test3
This is the first definition of the function name
This is a repeat of the same function name
This is the end of the script
$
```

The original definition of the `func1` function works fine, but after the second definition of the `func1` function, any subsequent uses of the function use the second definition.

Returning a Value

The bash shell treats functions like mini-scripts, complete with an exit status (see Chapter 10). There are three different ways you can generate an exit status for your functions.

The Default Exit Status

By default, the exit status of a function is the exit status returned by the last command in the function. After the function executes, you use the standard `$?` variable to determine the exit status of the function:

```
$ cat test4
#!/bin/bash
# testing the exit status of a function
```

```

func1() {
    echo "trying to display a non-existent file"
    ls -l badfile
}

echo "testing the function: "
func1
echo "The exit status is: $?"
$
$ ./test4
testing the function:
trying to display a non-existent file
ls: badfile: No such file or directory
The exit status is: 1
$

```

The exit status of the function is 1 because the last command in the function failed. However, you have no way of knowing if any of the other commands in the function completed successfully or not. Take a look at this example:

```

$ cat test4b
#!/bin/bash
# testing the exit status of a function

func1() {
    ls -l badfile
    echo "This was a test of a bad command"
}

echo "testing the function:"
func1
echo "The exit status is: $?"
$
$ ./test4b
testing the function:
ls: badfile: No such file or directory
This was a test of a bad command
The exit status is: 0
$

```

This time, because the function ended with an `echo` statement that completed successfully, the exit status of the function is 0, even though one of the commands in the function failed. Using the default exit status of a function can

be a dangerous practice. Fortunately, there are a couple of other solutions for us.

Using the return Command

The bash shell uses the `return` command to exit a function with a specific exit status. The `return` command allows you to specify a single integer value to define the function exit status, providing an easy way for you to programmatically set the exit status of your function:

```
$ cat test5
#!/bin/bash
# using the return command in a function

function dbl {
    read -p "Enter a value: " value
    echo "doubling the value"
    return ${ value * 2 }
}

dbl
echo "The new value is $?"
$
```

The `dbl` function doubles the value contained in the `$value` variable provided by the user input. It then returns the result using the `return` command, which the script displays using the `$?` variable.

You must be careful, however, when using this technique to return a value from a function. Keep the following two tips in mind to avoid problems:

- Remember to retrieve the return value as soon as the function completes.
- Remember that an exit status must be in the range of

0 to 255.

If you execute any other commands before retrieving the value of the function, using the `$?` variable, the return value from the function will be lost. Remember that the `$?` variable returns the exit status of the last executed command.

The second problem defines a limitation for using this return value technique. Because an exit status must be less than 256, the result of your function must produce an integer value less than 256. Any value over that returns an error value:

```
$ ./test5
Enter a value: 200
doubling the value
The new value is 1
$
```

You cannot use this return value technique if you need to return either larger integer values or a string value. Instead, you'll need to use another method, demonstrated in the next section.

Using Function Output

Just as you can capture the output of a command to a shell variable, you can also capture the output of a function to a shell variable. You can use this technique to retrieve any type of output from a function to assign to a variable:

```
result=`dbl`
```

This command assigns the output of the `dbl` function to the `$result` shell variable. The following is an example of using this method in a script:

```
$ cat test5b
#!/bin/bash
# using the echo to return a value

function dbl {
    read -p "Enter a value: " value
    echo ${ $value * 2 }
}

result=`dbl`
echo "The new value is $result"
$
$ ./test5b
Enter a value: 200
The new value is 400
$
$ ./test5b
Enter a value: 1000
The new value is 2000
$
```

The new function now uses an `echo` statement to display the result of the calculation. The script just captures the output of the `dbl` function instead of looking at the exit status for the answer.

There's a subtle trick that this example demonstrates. You'll notice that the `dbl` function really outputs two messages. The `read` command outputs a short message querying the user for the value. The bash shell script is smart enough to not consider this as part of the `STDOUT` output and ignores it. If you had used an `echo` statement to produce this query message to the user, it would have been captured by the shell variable as well as the output value.

Note

Using this technique, you can also return floating point and string values, making this an extremely

Using Variables in Functions

You might have noticed in the `test5` example in the previous section that we used a variable called `$value` within the function to hold the value that it processed. When you use variables in your functions, you need to be somewhat careful about how you define and handle them. This is a common cause of problems in shell scripts. This section goes over a few techniques for handling variables both inside and outside your shell script functions.

Passing Parameters to a Function

As we mentioned earlier in the “Returning a Value” section, the `bash` shell treats functions just like mini-scripts. This means that you can pass parameters to a function just like a regular script (see Chapter 13).

Functions can use the standard parameter environment variables to represent any parameters passed to the function on the command line. For example, the name of the function is defined in the `$0` variable, and any parameters on the function command line are defined using the variables `$1`, `$2`, and so on. You can also use the special variable `$#` to determine the number of parameters passed to the function.

When specifying the function in your script, you must provide the parameters on the same command line as the function, like this:

```
func1 $value1 10
```

The function can then retrieve the parameter values using the parameter environment variables. Here's an example of using this method to pass values to a function:

```
$ cat test6
#!/bin/bash
# passing parameters to a function

function addem {
    if [ $# -eq 0 ] || [ $# -gt 2 ]
    then
        echo -1
    elif [ $# -eq 1 ]
    then
        echo ${ $1 + $1 }
    else
        echo ${ $1 + $2 }
    fi
}

echo -n "Adding 10 and 15: "
value=`addem 10 15`
echo $value
echo -n "Let's try adding just one number: "
value=`addem 10`
echo $value
echo -n "Now trying adding no numbers: "
value=`addem`
echo $value
echo -n "Finally, try adding three numbers: "
value=`addem 10 15 20`
echo $value
$
$ ./test6
Adding 10 and 15: 25
Let's try adding just one number: 20
Now trying adding no numbers: -1
Finally, try adding three numbers: -1
$
```

The `addem` function in the `test6` script first checks the

number of parameters passed to it by the script. If there aren't any parameters, or if there are more than two parameters, `addem` returns a value of `-1`. If there's just one parameter, `addem` adds the parameter to itself for the result. If there are two parameters, `addem` adds them together for the result.

Because the function uses the special parameter environment variables for its own parameter values, it can't directly access the script parameter values from the command line of the script. The following example will fail:

```
$ cat badtest1
#!/bin/bash
# trying to access script parameters inside a function

function badfuncl {
    echo ${ $1 * $2 }
}

if [ $# -eq 2 ]
then
    value='badfuncl `
    echo "The result is $value"
else
    echo "Usage: badtest1 a b"
fi
$
$ ./badtest1
Usage: badtest1 a b
$ ./badtest1 10 15
./badtest1: * : syntax error: operand expected (error token is "**
")
The result is
$
```

Even though the function uses the `$1` and `$2` variables, they aren't the same `$1` and `$2` variables available in the main part of the script. Instead, if you want to use those values in your function, you'll have to manually pass them when you call the function:

```
$ cat test7
#!/bin/bash
# trying to access script parameters inside a function

function func7 {
    echo ${ $1 * $2 }
}

if [ $# -eq 2 ]
then
    value=`func7 $1 $2`
    echo "The result is $value"
else
    echo "Usage: badtest1 a b"
fi
$
$ ./test7
Usage: badtest1 a b
$ ./test7 10 15
The result is 150
$
```

By passing the `$1` and `$2` variables to the function, they become available for the function to use, just like any other parameter.

Handling Variables in a Function

One thing that causes problems for shell script programmers is the *scope* of a variable. The scope is where the variable is visible. Variables defined in functions can have a different scope than regular variables. That is, they can be hidden from the rest of the script.

Functions use two types of variables:

- Global
- Local

The following sections describe how to use both types of

variables in your functions.

Global Variables

Global variables are variables that are valid anywhere within the shell script. If you define a global variable in the main section of a script, you can retrieve its value inside a function. Likewise, if you define a global variable inside a function, you can retrieve its value in the main section of the script.

By default, any variables you define in the script are global variables. Variables defined outside of a function can be accessed within the function just fine:

```
$ cat test8
#!/bin/bash
# using a global variable to pass a value

function dbl {
    value=$(( $value * 2 ))
}

read -p "Enter a value: " value
dbl
echo "The new value is: $value"
$
$ ./test8
Enter a value: 450
The new value is: 900
$
```

The `$value` variable is defined outside of the function, and assigned a value outside of the function. When the `dbl` function is called, the variable and its value are still valid inside the function. When the variable is assigned a new value inside the function, that new value is still valid when the script references the variable.

This can be a dangerous practice, however, especially if you intend to use your functions in different shell scripts. It requires that you know exactly what variables are used in the function, including any variables used to calculate values not returned to the script. Here's an example of how things can go bad:

```
$ cat badtest2
#!/bin/bash
# demonstrating a bad use of variables

function func1 {
    temp=$(( $value + 5 ))
    result=$(( $temp * 2 ))
}

temp=4
value=6

func1
echo "The result is $result"
if [ $temp -gt $value ]
then
    echo "temp is larger"
else
    echo "temp is smaller"
fi
$
$ ./badtest2
The result is 22
temp is larger
$
```

Because the `$temp` variable was used in the function, its value is compromised in the script, producing a result that you may not have intended. There's an easy way to solve this problem in your functions, as shown in the next section.

Local Variables

Instead of using global variables in functions, any variables

that the function uses internally can be declared as local variables. To do that, just use the `local` keyword in front of the variable declaration:

```
local temp
```

You can also use the `local` keyword in an assignment statement while assigning a value to the variable:

```
local temp=$(( $value + 5 )
```

The `local` keyword ensures that the variable is limited to only within the function. If a variable with the same name appears outside the function in the script, the shell keeps the two variable values separate. Now you can easily keep your function variables separate from your script variables and share only the ones you want to share:

```
$ cat test9
#!/bin/bash
# demonstrating the local keyword
```

```
function func1 {
    local temp=$(( $value + 5 )
    result=$(( $temp * 2 )
}
```

```
temp=4
value=6
```

```
func1
echo "The result is $result"
if [ $temp -gt $value ]
then
    echo "temp is larger"
else
    echo "temp is smaller"
fi
$
$ ./test9
The result is 22
temp is smaller
$
```

Now when you use the `$temp` variable within the `func1` function, it doesn't affect the value assigned to the `$temp` variable in the main script.

Array Variables and Functions

Chapter 5 discussed an advanced way of allowing a single variable to hold multiple values by using arrays. Using array variable values with functions is a little tricky, and there are some special considerations. This section describes a technique that allows you to do that.

Passing Arrays to Functions

The art of passing an array variable to a script function can be confusing. If you try to pass the array variable as a single parameter, it won't work:

```
$ cat badtest3
#!/bin/bash
# trying to pass an array variable

function testit {
    echo "The parameters are: $@"
    thisarray=$1
    echo "The received array is ${thisarray[*]}"
}

myarray=(1 2 3 4 5)
echo "The original array is: ${myarray[*]}"
testit $myarray

$
$ ./badtest3
The original array is: 1 2 3 4 5
The parameters are: 1
./badtest3: thisarray[*]: bad array subscript
The received array is
```

If you try using the array variable as a function parameter, the function only picks up the first value of the array variable.

To solve this problem, you must disassemble the array variable into its individual values and then use the values as function parameters. Inside the function, you can reassemble all of the parameters into a new array variable. The following is an example of doing this:

```
$ cat test10
#!/bin/bash
# array variable to function test

function testit {
    local newarray
    newarray=('echo "$@"')
    echo "The new array value is: ${newarray[*]}"
}

myarray=(1 2 3 4 5)
echo "The original array is ${myarray[*]}"
testit ${myarray[*]}
$
$ ./test10
The original array is 1 2 3 4 5
The new array value is: 1 2 3 4 5
$
```

The script uses the *\$myarray* variable to hold all of the individual array values to place them all on the command line for the function. The function then rebuilds the array variable from the command line parameters. Once inside the function, the array can be used just like any other array:

```
$ cat test11
#!/bin/bash
# adding values in an array

function addarray {
    local sum=0
```

```

local newarray
newarray=('echo "$@"')
for value in ${newarray[*]}
do
    sum=$(( $sum + $value )) done
echo $sum
}

myarray=(1 2 3 4 5)
echo "The original array is: ${myarray[*]}"
arg1='echo ${myarray[*]}'
result='addarray $arg1'
echo "The result is $result"
$
$ ./test11
The original array is: 1 2 3 4 5
The result is 15
$

```

The `addarray` function iterates through the array values, adding them together. You can put any number of values in the `myarray` array variable, and the `addarray` function will add them.

Returning Arrays from Functions

Passing an array variable from a function back to the shell script uses a similar technique. The function uses an `echo` statement to output the individual array values in the proper order, and then the script must reassemble them into a new array variable:

```

$ cat test12
#!/bin/bash
# returning an array value

function arraydbl {
    local origarray
    local newarray
    local elements
    local i
    origarray=('echo "$@"')
    newarray=('echo "$@"')

```

```

elements=${# - 1 }
for (( i = 0; i <= $elements; i++ ))
{
    newarray[$i]=${origarray[$i]} * 2 }
}
echo ${newarray[*]}
}

myarray=(1 2 3 4 5)
echo "The original array is: ${myarray[*]}"
arg1='echo ${myarray[*]}'
result=('arraydbl $arg1')
echo "The new array is: ${result[*]}"
$
$ ./test12
The original array is: 1 2 3 4 5
The new array is: 2 4 6 8 10

```

The script passes the array value, using the `$arg1` variable to the `arraydbl` function. The `arraydbl` function reassembles the array into a new array variable, and it makes a copy for the output array variable. It then iterates through the individual array variable values, doubles each value, and places it into the copy of the array variable in the function.

The `arraydbl` function then uses the `echo` statement to output the individual values of the array variable values. The script uses the output of the `arraydbl` function to reassemble a new array variable with the values.

Function Recursion

One feature that local function variables provide is *self-containment*. A self-contained function doesn't use any resources outside of the function, other than whatever variables the script passes to it in the command line.

This feature enables the function to be called *recursively*, which means that the function calls itself to reach an answer. Usually, a recursive function has a base value that it eventually iterates down to. Many advanced mathematical algorithms use recursion to reduce a complex equation down one level repeatedly, until they get to the level defined by the base value.

The classic example of a recursive algorithm is calculating factorials. A factorial of a number is the value of the preceding numbers multiplied with the number. Thus, to find the factorial of 5, you'd perform the following equation:

$$5! = 1 * 2 * 3 * 4 * 5 = 120$$

Using recursion, the equation is reduced down to the following format:

$$x! = x * (x-1)!$$

or in English, the factorial of x is equal to x times the factorial of $x-1$. This can be expressed in a simple recursive script:

```
function factorial {
    if [ $1 -eq 1 ]
    then
        echo 1
    else
        local temp=$(( $1 - 1 ))
        local result=`factorial $temp`
        echo $[ $result * $1 ]
    fi
}
```

The factorial function uses itself to calculate the value for the factorial:

```
$ cat test13
#!/bin/bash
# using recursion
```

```
function factorial {
    if [ $1 -eq 1 ]
    then
        echo 1
    else
        local temp=$(( $1 - 1 ))
        local result='factorial $temp'
        echo ${result} * $1
    fi
}

read -p "Enter value: " value
result='factorial $value'
echo "The factorial of $value is: $result"
$
$ ./test13
Enter value: 5
The factorial of 5 is: 120
$
```

Using the factorial function is easy. Having created a function like this, you may want to use it in other scripts. Next, we'll look at how to do that efficiently.

Creating a Library

It's easy to see how functions can help save typing in a single script, but how about if you just happen to use the same single code block between scripts? It's obviously challenging if you have to define the same function in each script, only to use it one time in each script.

There's a solution for that problem! The bash shell allows you to create a *library file* for your functions, and then reference that single library file in as many scripts as you need to.

The first step in the process is to create a common

library file that contains the functions you need in your scripts. Here's a simple library file called `myfuncs` that defines three simple functions:

```
$ cat myfuncs
# my script functions

function addem {
    echo $[ $1 + $2 ]
}

function multem {
    echo $[ $1 * $2 ]
}

function divem {
    if [ $2 -ne 0 ]
    then
        echo $[ $1 / $2 ]
    else
        echo -1
    fi
}
$
```

The next step is to include the `myfuncs` library file in your script files that want to use any of the functions. This is where things get tricky.

The problem is with the scope of shell functions. As with environment variables, shell functions are only valid for the shell session in which you define them. If you run the `myfuncs` shell script from your shell command line interface prompt, the shell creates a new shell and runs the script in that new shell. This will define the three functions for that shell, but when you try to run another script that uses those functions, they won't be available.

This applies to scripts as well. If you try to just run the library file as a regular script file, the functions won't appear

in your script:

```
$ cat badtest4
#!/bin/bash
# using a library file the wrong way
./myfuncs

result='addem 10 15'
echo "The result is $result"
$
$ ./badtest4
./badtest3: addem: command not found
The result is
$
```

The key to using function libraries is the `source` command. The `source` command executes commands within the current shell context instead of creating a new shell to execute them. You use the `source` command to run the library file script inside of your shell script. This makes the functions available to the script.

The `source` command has a shortcut alias, called the *dot operator*. To source the `myfuncs` library file in a shell script, all you need to do is add the following line:

```
. ./myfuncs
```

This example assumes that the `myfuncs` library file is located in the same directory as the shell script. If not, you'll need to use the appropriate path to access the file. Here's an example of creating a script that uses the `myfuncs` library file:

```
$ cat test14
#!/bin/bash
# using functions defined in a library file
. ./myfuncs

value1=10
value2=5
result1='addem $value1 $value2'
```

```
result2='multem $value1 $value2`  
result3='divem $value1 $value2`  
echo "The result of adding them is: $result1"  
echo "The result of multiplying them is: $result2"  
echo "The result of dividing them is: $result3"  
$  
$ ./test14  
The result of adding them is: 15  
The result of multiplying them is: 50  
The result of dividing them is: 2  
$
```

The script successfully uses the functions defined in the `myfuncs` library file.

Using Functions on the Command Line

You can use script functions to create some pretty complex operations. Sometimes it would be nice to be able to use these functions directly on the command line interface prompt.

Just as you can use a script function as a command in a shell script, you can also use a script function as a command in the command line interface. This is a nice feature because once you define the function in the shell, you can use it from any directory on the system; you don't have to worry about a script being in your `PATH` environment variable. The trick is to get the shell to recognize the function. There are a couple of ways to do that.

Creating Functions on the Command

Line

Because the shell interprets commands as you type them, you can define a function directly on the command line. There are two ways to do that.

The first method defines the function all on one line:

```
$ function divem { echo $[ $1 / $2 ]; }
$ divem 100 5
20
$
```

When you define the function on the command line, you must remember to include a semicolon at the end of each command, so the shell knows where to separate commands:

```
$ function doubleit { read -p "Enter value: " value; echo $[
  $ value * 2 ]; }
$
$ doubleit
Enter value: 20
40
$
```

The other method is to use multiple lines to define the function. When you do that, the bash shell uses the secondary prompt to prompt you for more commands. Using this method, you don't need to place a semicolon at the end of each command; just press the Enter key:

```
$ function multem {
> echo $[ $1 * $2 ]
> }
$ multem 2 5
10
$
```

When you use the brace at the end of the function, the shell knows that you're done defining the function.

Caution

Be extremely careful when creating functions on the command line. If you use a function with the same name as a built-in command or another command, the function will override the original command.

Defining Functions in the `.bashrc` File

The obvious downside to defining shell functions directly on the command line is that when you exit the shell, your function disappears. For complex functions, this can become a problem.

A much simpler method is to define the function in a place where it will be reloaded by the shell each time you start a new shell.

The best place to do that is the `.bashrc` file. The bash shell looks for this file in your home directory each time it starts, whether interactively or as the result of starting a new shell from within an existing shell.

Directly Defining Functions

You can define the functions directly in the `.bashrc` file in your home directory. Most Linux distributions already define some things in the `.bashrc` file, so be careful not to remove those items. Just add your functions to the bottom of the existing file. Here's an example of doing that:

```
$ cat .bashrc
```

```
# .bashrc

# Source global definitions
if [ -r /etc/bashrc ]; then
    . /etc/bashrc
fi

function addem {
    echo $[ $1 + $2 ]
}

$
```

The function won't take effect until the next time you start a new bash shell. After you do that, you can use the function anywhere on the system.

Sourcing Function Files

Just as in a shell script, you can use the `source` command (or its alias the dot operator) to add functions from an existing library file to your `.bashrc` script:

```
$ cat .bashrc
# .bashrc

# Source global definitions
if [ -r /etc/bashrc ]; then
    . /etc/bashrc
fi

. /home/rich/libraries/myfuncs

$
```

Make sure that you include the proper pathname to reference the library file for the bash shell to find. The next time you start a shell, all of the functions in your library are available at the command line interface:

```
$ addem 10 5
15
$ multem 10 5
50
$ divem 10 5
```

\$

Even better, the shell also passes any defined functions to child shell processes so your functions are automatically available for any shell scripts you run from your shell session. You can test this by writing a script that uses the functions without defining or sourcing them:

```
$ cat test15
#!/bin/bash
# using a function defined in the .bashrc file

value1=10
value2=5
result1='addem $value1 $value2'
result2='multem $value1 $value2'
result3='divem $value1 $value2'
echo "The result of adding them is: $result1"
echo "The result of multiplying them is: $result2"
echo "The result of dividing them is: $result3"
$
$ ./test15
The result of adding them is: 15
The result of multiplying them is: 50
The result of dividing them is: 2
$
```

Even without sourcing the library file, the functions worked perfectly in the shell script.

Summary

Shell script functions allow you to place script code that's repeated throughout the script in a single place. Instead of having to rewrite blocks of code, you can create a function containing the code block, and then just reference the function name in your script. The bash shell jumps to the

function code block whenever it sees the function name used in the script.

You can even create script functions that return values. This allows you to create functions that interact with the script, returning both numeric and character data. Script functions can return numeric data by using the exit status of the last command in the function, or using the `return` command. The `return` command allows you to programmatically set the exit status of your function to a specific value based on the results of the function.

Functions can also return values using the standard `echo` statement. You can capture the output data using the backtick character as you would any other shell command. This enables you to return any type of data from a function, including strings and floating-point numbers.

You can use shell variables within your functions, assigning values to variables and retrieving values from existing variables. This allows you to pass any type of data both into and out of a script function from the main script program. Functions also allow you to define local variables, which are accessible only from within the function code block. Local variables allow you to create self-contained functions, which don't interfere with any variables or processes used in the main shell script.

Functions can also call other functions, including themselves. When a function calls itself, it is called recursion. A recursive function often has a base value that is the terminal value of the function. The function continues

to call itself with a decreasing parameter value until the base value is reached.

If you use lots of functions in your shell scripts, you can create library files of script functions. The library files can be included in any shell script file by using the source command, or its alias, the dot operator. This is called sourcing the library file. The shell doesn't run the library file but makes the functions available within the shell that runs the script. You can use this same technique to create functions that you can use on the normal shell command line. You can either define functions directly on the command line or you can add them to your `.bashrc` file so that they are available for each new shell session you start. This is a handy way to create utilities that can be used no matter what your `PATH` environment variable is set to.

The next chapter discusses the use of text graphics in your scripts. In this day of modern graphical interfaces, sometimes a plain text interface just doesn't cut it. The bash shell provides some easy ways for you to incorporate simple graphics features in your scripts to help spice things up.

Chapter 17

Writing Scripts for Graphical Desktops

In This Chapter

- Creating text menus
- Building text window widgets
- Adding X Windows graphics

Over the years, shell scripts have acquired a reputation for being dull and boring. This doesn't have to be the case, however, if you plan on running your scripts in a graphical environment. There are plenty of ways to interact with your script user that don't rely on the `read` and `echo` statements. This chapter dives into a few different methods you can use to help add life to your interactive scripts so that they don't look so old-fashioned.

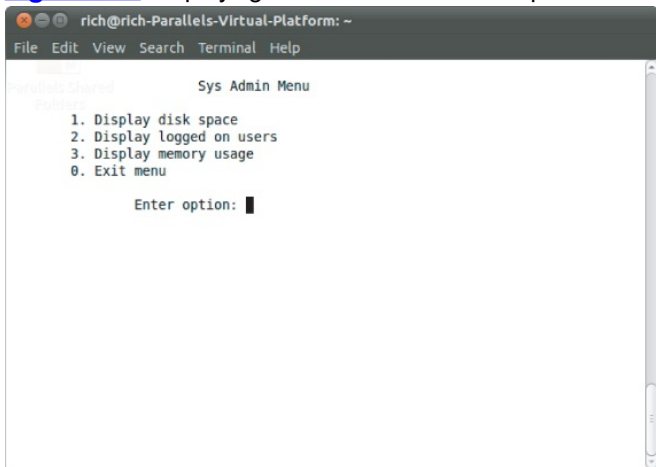
Creating Text Menus

The most common way to create an interactive shell script

is to utilize a menu. Offering your customers a choice of various options helps guide them through exactly what the script can and can't do.

Menu scripts usually clear the display area, and then show a list of options available. The customer can select an option by pressing an associated letter or number assigned to each option. [Figure 17.1](#) shows the layout of a sample menu.

[Figure 17.1](#) Displaying a menu from a shell script



The core of a shell script menu is the `case` command (see Chapter 11). The `case` command performs specific commands, depending on what character your customer

selects from the menu.

The following sections walk you through the steps you should follow to create a menu-based shell script.

Create the Menu Layout

The first step in creating a menu is, obviously, to determine what elements you want to appear in the menu and lay them out the way that you want them to appear.

Before creating the menu, it's usually a good idea to clear the monitor display. This enables you to display your menu in a clean environment without distracting text.

The `clear` command uses the terminfo data of your terminal session (see Chapter 2) to clear any text that appears on the monitor. After the `clear` command, you can use the `echo` command to display your menu elements.

By default, the `echo` command can only display printable text characters. When creating menu items, it's often helpful to use nonprintable items, such as the tab and newline characters. To include these characters in your `echo` command, you must use the `-e` option. Thus, the command:

```
echo -e "1.\tDisplay disk space"
```

results in the output line:

```
1.      Display disk space
```

This greatly helps in formatting the layout of the menu items. With just a few `echo` commands, you can create a reasonable-looking menu:

```
clear
echo
```

```
echo -e "\t\t\tSys Admin Menu\n"  
echo -e "\t1. Display disk space"  
echo -e "\t2. Display logged on users"  
echo -e "\t3. Display memory usage"  
echo -e "\t0. Exit menu\n\n"  
echo -en "\t\tEnter option: "
```

The `-en` option on the last line displays the line without adding the newline character at the end. This gives the menu a more professional look, as the cursor will stay at the end of the line waiting for the customer's input.

The last part of creating the menu is to retrieve the input from the customer. This is done using the `read` command (see Chapter 13). Because we expect only single-character input, the nice thing to do is to use the `-n` option in the `read` command to retrieve only one character. This allows the customer to enter a number without having to press the Enter key:

```
read -n 1 option
```

Next, you'll need to create your menu functions.

Create the Menu Functions

Shell script menu options are easier to create as a group of separate functions. This enables you to create a simple, concise `case` command that is easy to follow.

To do that, you need to create separate shell functions for each of your menu options. The first step in creating a menu shell script is to determine what functions you want your script to perform and lay them out as separate functions in your code.

It is common practice to create *stub functions* for functions that aren't implemented yet. A stub function is a function that doesn't contain any commands yet, or possibly just an `echo` statement indicating what should be there eventually:

```
function diskpace {  
    clear  
    echo "This is where the diskpace commands will go"  
}
```

This enables your menu to operate smoothly while you work on the individual functions. You don't have to code all of the functions for your menu to work. You'll notice that the function starts out with the `clear` command. This enables you to start the function on a clean monitor screen, without the menu showing.

One thing that helps out in the shell script menu is to create the menu layout itself as a function:

```
function menu {  
    clear  
    echo  
    echo -e "\t\t\tSys Admin Menu\n"  
    echo -e "\t1. Display disk space"  
    echo -e "\t2. Display logged on users"  
    echo -e "\t3. Display memory usage"  
    echo -e "\t0. Exit program\n\n"  
    echo -en "\t\tEnter option: "  
    read -n 1 option  
}
```

This enables you to easily redisplay the menu at any time just by calling the `menu` function.

Add the Menu Logic

Now that you have your menu layout and your functions, all

you need to do is create the programming logic to put the two together. As mentioned, this requires the `case` command.

The `case` command should call the appropriate function according to the character selection expected from the menu. It's always a good idea to use the default `case` command character (the asterisk) to catch any incorrect menu entries.

The following code illustrates the use of the `case` command in a typical menu:

```
menu
case $option in
0)
    break ;;
1)
    diskspace ;;
2)
    whoseon ;;
3)
    memusage ;;
*)
    clear
    echo "Sorry, wrong selection";;
esac
```

This code first uses the `menu` function to clear the monitor screen and display the menu. The `read` command in the `menu` function pauses until the customer hits a character on the keyboard. Once that's been done, the `case` command takes over. The `case` command calls the appropriate function based on the returned character. After the function completes, the `case` command exits.

Putting It All Together

Now that you've seen all of the parts that make up a shell script menu, let's put them together and see how they all interoperate. Here's an example of a full menu script:

```
$ cat menu1
#!/bin/bash
# simple script menu

function diskspace {
    clear
    df -k
}

function whoseon {
    clear
    who
}

function memusage {
    clear
    cat /proc/meminfo
}

function menu {
    clear
    echo
    echo -e "\t\t\tSys Admin Menu\n"
    echo -e "\t1. Display disk space"
    echo -e "\t2. Display logged on users"
    echo -e "\t3. Display memory usage"
    echo -e "\t0. Exit program\n\n"
    echo -en "\t\t\tEnter option: "
    read -n 1 option
}

while [ 1 ]
do
    menu
    case $option in
        0)
            break ;;
        1)
            diskspace ;;
        2)
            whoseon ;;
        3)
            memusage ;;
        *)
```

```
clear
echo "Sorry, wrong selection";;
esac
echo -en "\n\n\t\t\tHit any key to continue"
read -n 1 line
done
clear
$
```

This menu creates three functions to retrieve administrative information about the Linux system using common commands. It uses a `while` loop to continually loop through the menu until the customer selects option 0, which uses the `break` command to break out of the `while` loop.

You can use this same template to create any shell script menu interface. It provides a simple way to interact with your customers.

Using the select Command

You may have noticed that half the challenge of creating a text menu is just creating the menu layout and retrieving the answer that you enter. The bash shell provides a handy little utility for you that does all of this work automatically.

The `select` command allows you to create a menu from a single command line and then retrieve the entered answer and automatically process it. The format of the `select` command is as follows:

```
select variable in list
do
    commands
done
```

The `list` parameter is a space-separated list of text items that build the menu. The `select` command displays each

item in the list as a numbered option and then displays a special prompt, defined by the `PS3` environment variable, for the selection.

Here's a simple example of the `select` command in action:

```
$ cat smenu1
#!/bin/bash
# using select in the menu

function diskspace {
    clear
    df -k
}

function whoseon {
    clear
    who
}

function memusage {
    clear
    cat /proc/meminfo
}

PS3="Enter option: "
select option in "Display disk space" "Display logged on users"
"Display memory usage" "Exit program"
do
    case $option in
        "Exit program")
            break ;;
        "Display disk space")
            diskspace ;;
        "Display logged on users")
            whoseon ;;
        "Display memory usage")
            memusage ;;
        *)
            clear
            echo "Sorry, wrong selection";;
    esac
done
clear
$
```

When you run the program, it automatically produces the

following menu:

```
$ ./smenu1
1) Display disk space      3) Display memory usage
2) Display logged on users 4) Exit program
Enter option:
```

When you use the `select` command, remember that the result value stored in the variable is the entire text string and not the number associated with the menu item. The text string values are what you need to compare in your `case` statements.

Doing Windows

Using text menus is a step in the right direction, but there's still a lot missing in our interactive scripts, especially if we try to compare them to the graphical Windows world. Fortunately for us, there are some very resourceful people out in the open source world that have helped us out.

The *dialog* package is a nifty little tool originally created by Savio Lam, and currently maintained by Thomas E. Dickey. This package recreates standard Windows dialog boxes in a text environment using ANSI escape control codes. You can easily incorporate these dialog boxes in your shell scripts to interact with your script users. This section describes the *dialog* package and demonstrates how to use it in shell scripts.

Note

The `dialog` package isn't installed in all Linux distributions by default. If it's not installed by default, because of its popularity it's almost always included in the software repository. Check your specific Linux distribution documentation for how to load the `dialog` package. For the Ubuntu Linux distribution, the following is the command line command to install it:

```
sudo apt-get install dialog
```

That package installs the `dialog` package plus the required libraries for your system.

The `dialog` Package

The `dialog` command uses command line parameters to determine what type of Windows *widget* to produce. A widget is the `dialog` package term for a type of Windows element. The `dialog` package currently supports the types of widgets shown in [Table 17.1](#).

Table 17.1 The `dialog` Widgets

Widget	Description
<code>calendar</code>	Provides a calendar to select a date from
<code>checklist</code>	Displays multiple entries where each entry can be turned on or off
<code>form</code>	Allows you to build a form with labels and text fields to be filled out
<code>fselect</code>	Provides a file selection window to browse for a file
<code>gauge</code>	Displays a meter showing a percentage of completion
<code>infobox</code>	Displays a message without waiting for a response
<code>inputbox</code>	Displays a single text form box for text entry

inputmenu	Provides an editable menu
menu	Displays a list of selections to choose from
msgbox	Displays a message and requires the user to select an OK button
pause	Displays a meter showing the status of a specified pause period
passwordbox	Displays a single textbox that hides entered text
passwordform	Displays a form with labels and hidden text fields
radiolist	Provides a group of menu items where only one item can be selected
tailbox	Displays text from a file in a scroll window using the <code>tail</code> command
tailboxbg	Same as <code>tailbox</code> , but operates in background mode
textbox	Displays the contents of a file in a scroll window
timebox	Provides a window to select an hour, minute, and second
yesno	Provides a simple message with Yes and No buttons

As you can see from [Table 17.1](#), there are lots of different widgets to choose from. This can give your scripts a more professional look with very little effort.

To specify a specific widget on the command line, you need to use the double dash format:

```
dialog --widget parameters
```

where *widget* is the widget name as seen in [Table 17.1](#), and *parameters* defines the size of the widget window and any text required for the widget.

Each dialog widget provides output in two forms:

- Using `STDERR`

- Using the exit code status

The exit code status of the `dialog` command determines the button selected by the user. If an OK or Yes button is selected, the `dialog` command returns a 0 exit status. If a Cancel or No button is selected, the `dialog` command returns a 1 exit status. You can use the standard `$?` variable to determine which button was selected in the dialog widget.

If a widget returns any data, such as a menu selection, the `dialog` command sends the data to `STDERR`. You can use the standard bash shell technique of redirecting the `STDERR` output to another file or file descriptor:

```
dialog --inputbox "Enter your age:" 10 20 2>age.txt
```

This command redirects the text entered in the textbox to the `age.txt` file.

The following sections take a look at some examples of the more common dialog widgets you'll use in your shell scripts.

The msgbox Widget

The `msgbox` widget is the most common type of dialog box. It displays a simple message in a window and waits for the user to click an OK button before disappearing. The following format is required to use a `msgbox` widget:

```
dialog --msgbox text height width
```

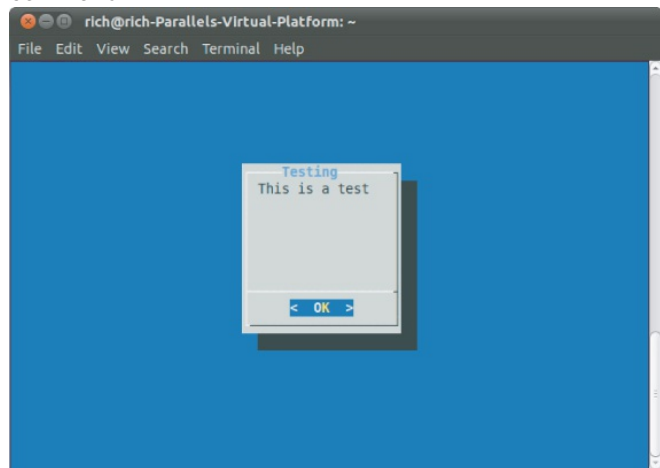
The `text` parameter is any string you want to place in the window. The `dialog` command will automatically wrap the

text to fit the size of the window you create, using the *height* and *width* parameters. If you want to place a title at the top of the window, you can also use the `--title` parameter, along with the text of the title. Here's an example of using the `msgbox` widget:

```
$ dialog --title Testing --msgbox "This is a test" 10 20
```

After entering this command, the message box will appear on the screen of the terminal emulator session you're using. [Figure 17.2](#) shows what this looks like.

Figure 17.2 Using the `msgbox` widget in the `dialog` command



If your terminal emulator supports the mouse, you can

click the OK button to close the dialog box. You can also use keyboard commands to simulate a click—just press the Enter key.

The yesno Widget

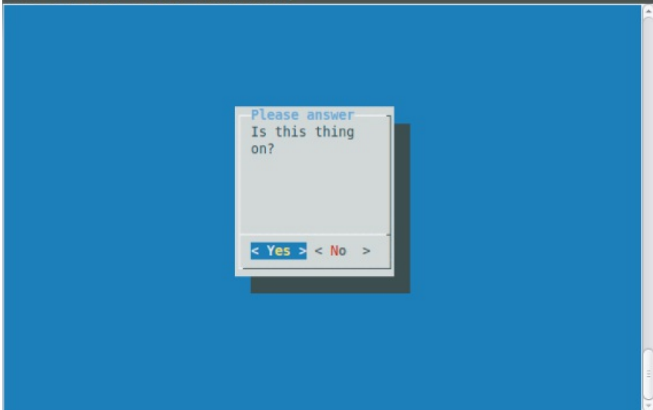
The `yesno` widget takes the `msgbox` widget one step further, allowing the user to answer a yes/no question displayed in the window. It produces two buttons at the bottom of the window, one for Yes, and another for No. The user can switch between buttons by using the mouse, the tab key, or the keyboard arrow keys. To select the button, the user can either press the spacebar or the Enter key.

Here's an example of using the `yesno` widget:

```
$ dialog --title "Please answer" --yesno "Is this thing on?" 10 20
$ echo $?
1
$
```

This produces the widget shown in [Figure 17.3](#).

[Figure 17.3](#) Using the `yesno` widget in the dialog command

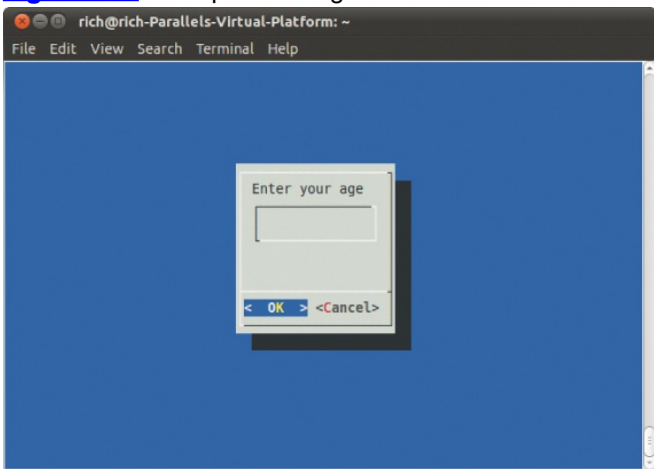


The exit status of the `dialog` command is set depending on which button the user selects. If the No button is selected, the exit status is 1, and if the Yes button is selected, the exit status is 0.

The inputbox Widget

The `inputbox` widget provides a simple textbox area for the user to enter a text string. The `dialog` command sends the value of the text string to `STDERR`. You must redirect that to retrieve the answer. [Figure 17.4](#) demonstrates what the `inputbox` widget looks like.

[Figure 17.4](#) The inputbox widget



As you can see in [Figure 17.4](#), the `inputbox` provides two buttons, OK and Cancel. If the Cancel button is selected, the exit status of the command is 1; otherwise, the exit status will be 0:

```
$ dialog --inputbox "Enter your age:" 10 20 2>age.txt
$ echo $?
0
$ cat age.txt
12$
```

You'll notice when you use the `cat` command to display the contents of the text file that there's no newline character after the value. This enables you to easily redirect the file contents to a variable in a shell script to extract the string

entered by the user.

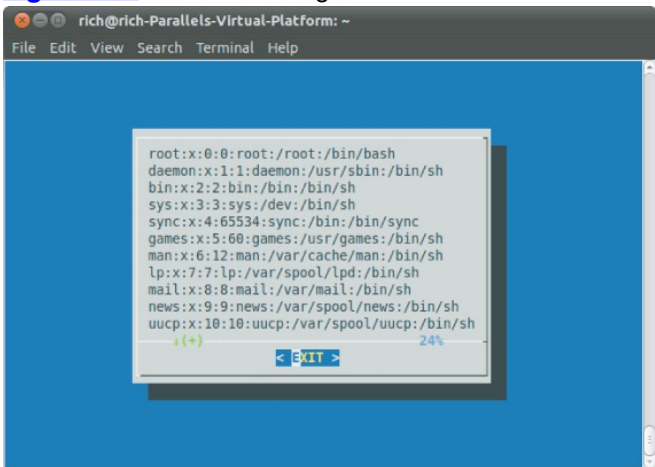
The *textbox* Widget

The `textbox` widget is a great way to display lots of information in a window. It produces a scrollable window containing the text from a file specified in the parameters:

```
$ dialog --textbox /etc/passwd 15 45
```

The contents of the `/etc/passwd` file are shown within the scrollable text window, as illustrated in [Figure 17.5](#).

Figure 17.5 The `textbox` widget



You can use the arrow keys to scroll left and right, as well

as up and down in the text file. The bottom line in the window shows the percent location within the file that you're viewing. The textbox contains only a single Exit button, which should be selected to exit the widget.

The menu Widget

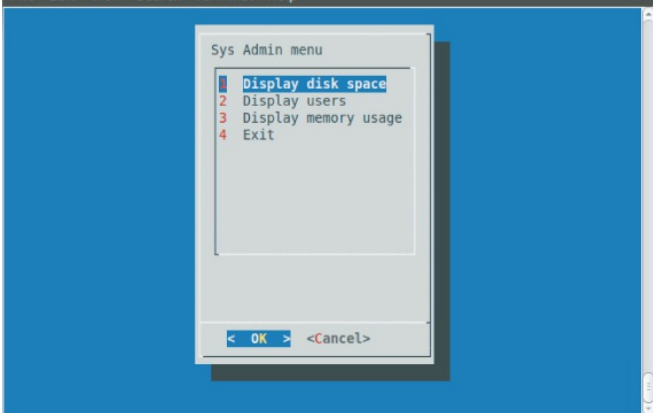
The `menu` widget allows you to create a window version of the text menu we created earlier in this chapter. All you need to do is provide a selection tag and the text for each item:

```
$ dialog --menu "Sys Admin Menu" 20 30 10 1 "Display disk space"  
2 "Display users" 3 "Display memory usage" 4 "Exit" 2> test.txt
```

The first parameter defines a title for the menu. The next two parameters define the height and width of the menu window, while the third parameter defines the number of menu items that appear in the window at one time. If there are more menu items, you can scroll through them using the arrow keys.

Following those parameters, you must add menu item pairs. The first element is the tag used to select the menu item. Each tag should be unique for each menu item and can be selected by hitting the appropriate key on the keyboard. The second element is the text used in the menu. [Figure 17.6](#) demonstrates the menu produced by the example command.

[Figure 17.6](#) The menu widget with menu items



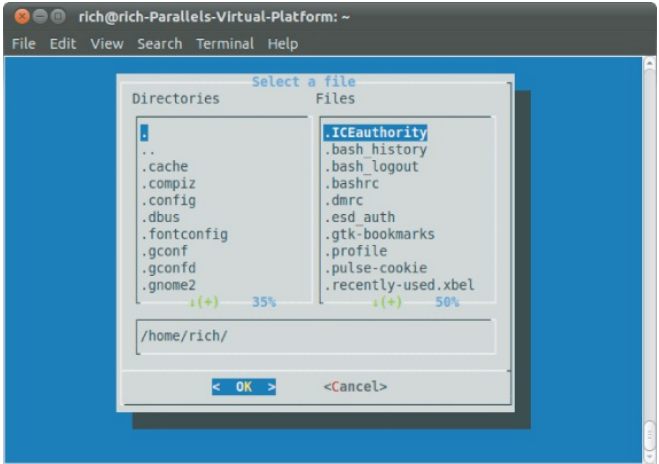
If the user selects a menu item by pressing the appropriate key for the tag, that menu item is highlighted but not selected. A selection isn't made until the OK button is selected by using either the mouse or the Enter key. The `dialog` command sends the selected menu item text to `STDERR`, which you can redirect as needed.

The fselect Widget

There are several fancy built-in widgets provided by the `dialog` command. The `fselect` widget is extremely handy when working with file names. Instead of forcing the user to type in a file name, you can use the `fselect` widget to browse

to the file location and select the file, as shown in [Figure 17.7](#).

Figure 17.7 The `fselect` widget



The `fselect` widget format looks like:

```
$ dialog --title "Select a file" --fselect $HOME/ 10 50 2>file.txt
```

The first parameter after the `fselect` option is the starting folder location used in the window. The `fselect` widget window consists of a directory listing on the left side, a file listing on the right side that shows all of the files in the selected directory, and a simple textbox that contains the currently selected file or directory. You can manually type a file name in the textbox, or use the directory and file listings

to select one.

The dialog Options

In addition to the standard widgets, you can customize lots of different options in the dialog command. You've already seen the `--title` parameter in action. This allows you to set a title for the widget that appears at the top of the window.

There are lots of other options available that allow you to completely customize both the appearance and the behavior of your windows. [Table 17.2](#) shows the options available for the `dialog` command.

Table 17.2 The dialog Command Options

Option	Description
<code>--add-widget</code>	Proceed to the next dialog unless Esc or the Cancel button has been pressed.
<code>--aspect ratio</code>	Specify the width/height aspect ratio of the window.
<code>--backtitle title</code>	Specify a title to display on the background, at the top of the screen.
<code>--begin x y</code>	Specify the starting location of the top-left corner of the window.
<code>--cancel-label label</code>	Specify an alternative label for the Cancel button.
<code>--clear</code>	Clear the display using the default dialog background color.
<code>--colors</code>	Embed ANSI color codes in dialog text.
<code>--cr-wrap</code>	Allow newline characters in dialog text and force a line wrap.
<code>--create-rc file</code>	Dump a sample configuration file to the specified file.

<code>--defaultno</code>	Make the default of a yes/no dialog No.
<code>--default-item string</code>	Set the default item in a checklist, form, or menu dialog.
<code>--exit-label label</code>	Specify an alternative label for the Exit button.
<code>--extra-button</code>	Display an extra button between the OK and Cancel buttons.
<code>--extra-label label</code>	Specify an alternative label for the Extra button.
<code>--help</code>	Display the dialog command help message.
<code>--help-button</code>	Display a Help button after the OK and Cancel buttons.
<code>--help-label label</code>	Specify an alternative label for the Help button.
<code>--help-status</code>	Write the checklist, radiolist, or form information after the help information in the Help button was selected.
<code>--ignore</code>	Ignore options that dialog does not recognize.
<code>--input-fd fd</code>	Specify an alternative file descriptor, other than STDIN.
<code>--insecure</code>	Change the password widget to display asterisks when typing.
<code>--item-help</code>	Add a help column at the bottom of the screen for each tag in a checklist, radiolist, or menu for the tag item.
<code>--keep-window</code>	Don't clear old widgets from the screen.
<code>--max-input size</code>	Specify a maximum string size for the input. The default is 2048.
<code>--nocancel</code>	Suppress the Cancel button.
<code>--no-collapse</code>	Don't convert tabs to spaces in dialog text.
<code>--no-kill</code>	Place the tailboxbg dialog in background and disable SIGHUP for the process.
<code>--no-label label</code>	Specify an alternative label for the No button.

<code>--no-shadow</code>	Don't display shadows for dialog windows.
<code>--ok-label label</code>	Specify an alternative label for the OK button.
<code>--output-fd fd</code>	Specify an alternative output file descriptor other than <code>STDERR</code> .
<code>--print-maxsize</code>	Print the maximum size of dialog windows allowed to the output.
<code>--print-size</code>	Print the size of each dialog window to the output.
<code>--print-version</code>	Print the dialog version to output.
<code>--separate-output</code>	Output the result of a checklist widget one line at a time with no quoting.
<code>--separator string</code>	Specify a string that separates the output for each widget.
<code>--separate-widget widget string</code>	Specify a string that separates the output for each widget.
<code>--shadow</code>	Draw a shadow to the right and bottom of each window.
<code>--single-quoted</code>	Use single quoting if needed for the checklist output.
<code>--sleep sec</code>	Delay for the specified number of seconds after processing the dialog window.
<code>--stderr</code>	Send output to <code>STDERR</code> (this is the default behavior).
<code>--stdout</code>	Send output to <code>STDOUT</code> .
<code>--tab-correct</code>	Convert tabs to spaces.
<code>--tab-len n</code>	Specify the number of spaces a tab character uses (the default is 8).
<code>--timeout sec</code>	Specify the number of seconds before exiting with an error code if no user input.
<code>--title title</code>	Specify the title of the dialog window.

<code>--trim</code>	Remove leading spaces and newline characters from dialog text.
<code>--visit-items</code>	Modify the tab stops in the dialog window to include the list of items.
<code>--yes-label</code> <code>label</code>	Specify an alternative label for the Yes button.

The `--backtitle` option is a handy way to create a common title for your menu through the script. If you specify it for each dialog window, it will persist throughout your application, creating a professional look to your script.

As you can tell from [Table 17.2](#), you can overwrite any of the button labels in your dialog window. This feature allows you to create just about any window situation you need.

Using the dialog Command in a Script

Using the `dialog` command in your scripts is a snap. There are just two things you must remember:

- Check the exit status of the dialog command if there's a Cancel or No button available.
- Redirect `STDERR` to retrieve the output value.

If you follow these two rules, you'll have a professional-looking interactive script in no time. Here's an example using dialog widgets to reproduce the system admin menu we created earlier in the chapter:

```
$ cat menu3
#!/bin/bash
# using dialog to create a menu

temp='mktemp -t test.XXXXXX'
temp2='mktemp -t test2.XXXXXX'
```

```

function diskpace {
    df -k > $temp
    dialog --textbox $temp 20 60
}

function whoseon {
    who > $temp
    dialog --textbox $temp 20 50
}

function memusage {
    cat /proc/meminfo > $temp
    dialog --textbox $temp 20 50
}

while [ 1 ]
do
    dialog --menu "Sys Admin Menu" 20 30 10 1 "Display disk space" 2
    "Display users" 3 "Display memory usage" 0 "Exit" 2> $temp2
    if [ $? -eq 1 ]
    then
        break
    fi

    selection='cat $temp2'

    case $selection in
    1)
        diskpace ;;
    2)
        whoseon ;;
    3)
        memusage ;;
    0)
        break ;;
    *)
        dialog --msgbox "Sorry, invalid selection" 10 30
    esac
done
rm -f $temp 2> /dev/null
rm -f $temp2 2> /dev/null
$

```

The script uses the `while` loop with a constant true value to create an endless loop displaying the menu dialog. This means that, after every function, the script returns to

displaying the menu.

The `menu` dialog includes a Cancel button, so the script checks the exit status of the dialog command in case the user presses the Cancel button to exit. Because it's in a `while` loop, exiting is as easy as using the `break` command to jump out of the `while` loop.

The script uses the `mktemp` command to create two temporary files for holding data for the `dialog` commands. The first one, `$temp`, is used to hold the output of the `df` and `meminfo` commands so that they can be displayed in the `textbox` dialog (see [Figure 17.8](#)). The second temporary file, `$temp2`, is used to hold the selection value from the main menu dialog.

[Figure 17.8](#) The `meminfo` command output displayed using the `textbox` dialog option

```
MemTotal:      1025296 kB
MemFree:       467776 kB
Buffers:       62560 kB
Cached:        304636 kB
SwapCached:    0 kB
Active:        245632 kB
Inactive:      255276 kB
Active(anon):  134024 kB
Inactive(anon): 3844 kB
Active(file):  111608 kB
Inactive(file): 251432 kB
Unevictable:   12 kB
Mlocked:      12 kB
HighTotal:     139208 kB
HighFree:      244 kB
LowTotal:      886088 kB
```

+ (+)

35%

< EXIT >

Now this is starting to look like a real application you can show off to people!

Getting Graphic

If you're looking for even more graphics for your interactive scripts, you can go one step further. Both the KDE and GNOME desktop environments (see Chapter 1) have expanded on the `dialog` command idea and include commands that produce X Windows graphical widgets for their respective environments.

This section describes the `kdiallog` and `zenity` packages,

which provide graphical window widgets for the KDE and GNOME desktops, respectively.

The KDE Environment

The KDE graphical environment includes the `kdIALOG` package by default. The `kdIALOG` package uses the `kdIALOG` command to generate standard windows, similar to the `dialog`-style widgets, within your KDE desktop. However, instead of having the clunky feel to them, these windows blend right in with the rest of your KDE application windows! This allows you to produce Windows-quality user interfaces directly from your shell scripts!

Note

Just because your Linux distribution uses the KDE desktop doesn't necessarily mean it has the `kdIALOG` package installed by default. You may need to manually install it from the distribution repository.

kdIALOG Widgets

Just like the `dialog` command, the `kdIALOG` command uses command line options to specify what type of window widget to use. The following is the format of the `kdIALOG` command:

```
kdIALOG display-options window-options arguments
```

The `window-options` options allow you to specify what type of

window widget to use. The available options are shown in [Table 17.3](#).

Table 17.3 kdialog Window Options

Option	Description
--checklist <i>title</i> [<i>tag item status</i>]	A checklist menu, with status specifying if the item is checked or not.
--error <i>text</i>	Error message box.
--inputbox <i>text</i> [<i>init</i>]	Input textbox. You can specify the default value using the <i>init</i> value.
--menu <i>title</i> [<i>tag item</i>]	Menu selection box title and a list of items identified by a tag.
--msgbox <i>text</i>	Simple message box with specified text.
--password <i>text</i>	Password input textbox that hides user input.
--radiolist <i>title</i> [<i>tag item status</i>]	A radiolist menu, with status specifying if the item is selected or not.
--separate-output	Returns items on separate lines for checklist and radiolist menus.
--sorry <i>text</i>	Sorry message box.
--textbox <i>file</i> [<i>width</i>] [<i>height</i>]	Textbox displaying the contents of <i>file</i> , alternatively specified by <i>width</i> and <i>height</i> .
--title <i>title</i>	Specifies a title for the TitleBar area of the dialog window.
--warningyesno <i>text</i>	Warning message box with Yes and No buttons.
--warningcontinuecancel <i>text</i>	Warning message box with Continue and Cancel buttons.
--warningyesnocancel <i>text</i>	Warning message box with Yes, No, and Cancel buttons.
--yesno <i>text</i>	Question box with Yes and No buttons.

As you can see from [Table 17.3](#), all of the standard window dialog box types are represented. However, when you use a `kdialog` window widget, it appears as a separate window in the KDE desktop, not inside the terminal emulator session!

The `checklist` and `radiolist` widgets allow you to define individual items in the lists and whether they are selected by default:

```
$kdialog --checklist "Items I need" 1 "Toothbrush" on 2 "Toothpaste"  
off 3 "Hair brush" on 4 "Deoderant" off 5 "Slippers" off
```

The resulting checklist window is shown in [Figure 17.9](#).

Figure 17.9 A `kdialog` checklist dialog window



The items specified as “on” are highlighted in the checklist. To select or deselect an item in the checklist, just click on it. If you select the OK button, the `kdiallog` will send the tag values to `STDOUT`:

```
"1" "3" "5"  
$
```

When you hit the Enter key, the `kdiallog` box appears with the selections. When you click the OK or Cancel buttons, the `kdiallog` command returns each tag as a string value to `STDOUT` (these are the “1”, “3”, and “5” values you see in the output). Your script must be able to parse the resulting values and match them with the original values.

Using *kdiallog*

You can use the `kdiallog` window widgets in your shell scripts similarly to how you use the `dialog` widgets. The big difference is that the `kdiallog` window widgets output values using `STDOUT` instead of `STDERR`.

Here's a script that converts the sys admin menu created earlier into a KDE application:

```
$ cat menu4  
#!/bin/bash  
# using kdiallog to create a menu  
  
temp='mktemp -t temp.XXXXXX'  
temp2='mktemp -t temp2.XXXXXX'  
  
function diskpace {  
    df -k > $temp  
    kdiallog --textbox $temp 1000 10  
}  
  
function whoseon {  
    who > $temp
```



```

        kdialog --textbox $temp 500 10
    }

function memusage {
    cat /proc/meminfo > $temp
    kdialog --textbox $temp 300 500
}

while [ 1 ]
do
kdialog --menu "Sys Admin Menu" "1" "Display diskspace" "2" "Display
users" "3" "Display memory usage" "0" "Exit" > $temp2
if [ $? -eq 1 ]
then
    break
fi

selection='cat $temp2'

case $selection in
1)
    diskspace ;;
2)
    whoseon ;;
3)
    memusage ;;
0)
    break ;;
*)
    kdialog --msgbox "Sorry, invalid selection"
esac
done
$

```

There isn't much difference in the script from using the `kdialog` command and the `dialog` command. The resulting main menu generated is shown in [Figure 17.10](#).

[Figure 17.10](#) The sys admin menu script using `kdialog`



Now your simple shell script looks just like a real KDE application! There's no limit to what you can do with your interactive scripts now.

The GNOME Environment

The GNOME graphical environment supports two popular packages that can generate standard windows:

- `gdialog`
- `zenity`

By far, `zenity` is the most commonly available package found in most GNOME desktop Linux distributions (it's installed by default in both Ubuntu and Fedora). This

section describes the features of zenity and demonstrates how to use it in your shell scripts.

zenity Widgets

As you would expect, zenity allows you to create different windows widgets by using command line options. [Table 17.4](#) shows the different widgets that zenity can produce.

Table 17.4 The zenity Windows Widgets

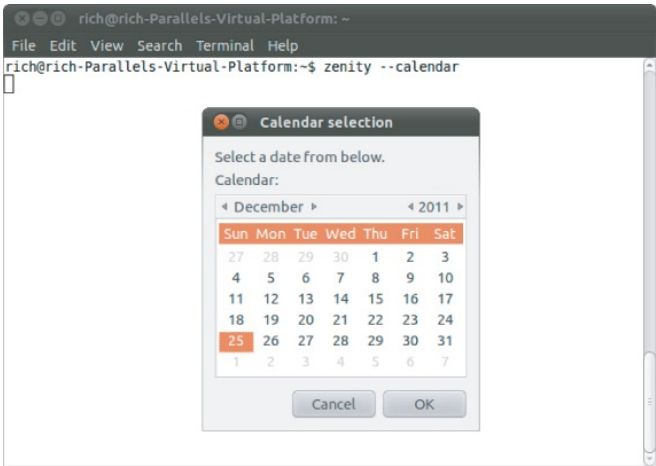
Option	Description
--calendar	Display a full month calendar.
--entry	Display a text entry dialog window.
--error	Display an error message dialog window.
--file-selection	Display a full pathname and file name dialog window.
--info	Display an informational dialog window.
--list	Display a checklist or radiolist dialog window.
--notification	Display a notification icon.
--progress	Display a progress bar dialog window.
--question	Display a yes/no question dialog window.
--scale	Display a scale dialog window.
--text-info	Display a textbox containing text.
--warning	Display a warning dialog window.

The `zenity` command line program works somewhat differently than the `kdiallog` and `dialog` programs. Many of the widget types are defined using additional options on the command line, instead of including them as arguments to

an option.

The `zenity` command does offer some pretty cool advanced dialog windows. The `calendar` option produces a full month calendar, as shown in [Figure 17.11](#).

Figure 17.11 The `zenity` calendar dialog window

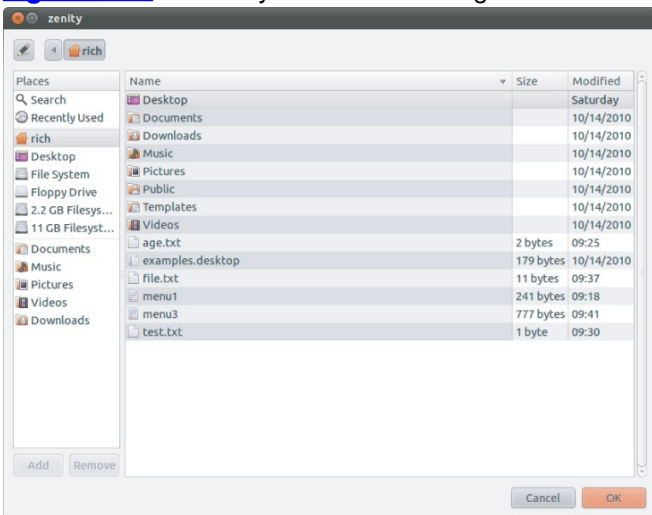


When you select a date from the calendar, the `zenity` command returns the value to `STDOUT`, just like `kdiallog`:

```
$ zenity --calendar
12/25/2011
$
```

Another pretty cool window in `zenity` is the file selection option, shown in [Figure 17.12](#).

Figure 17.12 The zenity file selection dialog window



You can use the dialog window to browse to any directory location on the system (as long as you have the privileges to view the directory) and select a file. When you select a file, the `zenity` command returns the full file and pathname:

```
$ zenity --file-selection
/home/ubuntu/menu5
$
```

With tools like that at your disposal, the sky's the limit with your shell script creations!

Using zenity in Scripts

As you would expect, zenity performs well in shell scripts. Unfortunately zenity chose not to follow the option convention used in `dialog` and `kdiallog`, so converting any existing interactive scripts to zenity may prove challenging.

In converting the sys admin menu from `kdiallog` to `zenity`, we found that we had to do quite a bit of manipulation of the widget definitions:

```
$cat menu5
#!/bin/bash
# using zenity to create a menu

temp='mktemp -t temp.XXXXXX'
temp2='mktemp -t temp2.XXXXXX'

function diskpace {
    df -k > $temp
    zenity --text-info --title "Disk space" --filename=$temp
    --width 750 --height 10
}

function whoseon {
    who > $temp
    zenity --text-info --title "Logged in users" --filename=$temp
    --width 500 --height 10
}

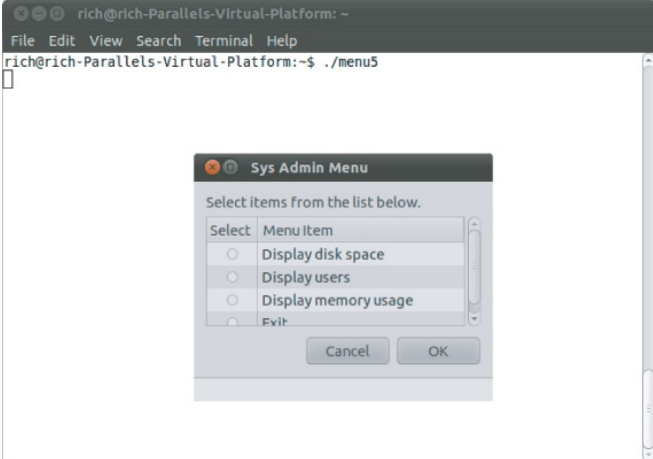
function memusage {
    cat /proc/meminfo > $temp
    zenity --text-info --title "Memory usage" --filename=$temp
    --width 300 --height 500
}

while [ 1 ]
do
    zenity --list --radiolist --title "Sys Admin Menu" --column "Select"
--column "Menu Item" FALSE "Display diskpace" FALSE "Display users"
FALSE "Display memory usage" FALSE "Exit" > $temp2
    if [ $? -eq 1 ]
    then
        break
    fi
}
```

```
selection='cat $temp2`  
case $selection in  
  "Display disk space")  
    diskspace ;;  
  "Display users")  
    whoseon ;;  
  "Display memory usage")  
    memusage ;;  
Exit)  
  break ;;  
*)  
  zenity --info "Sorry, invalid selection"  
esac  
done  
$
```

Because zenity doesn't support the menu dialog window, we used a radiolist type window for the main menu, as shown in [Figure 17.13](#).

[Figure 17.13](#) The sys admin menu using zenity



The radiolist uses two columns, each with a column heading. The first column includes the radio buttons to select. The second column is the item text. The radiolist also doesn't use tags for the items. When you select an item, the full text of the item is returned to `STDOUT`. This makes life a little more interesting for the `case` command. You must use the full text from the items in the `case` options. If there are any spaces in the text, you need to use quotation marks around the text.

Using the zenity package, you can add a Windows feel to your interactive shell scripts in the GNOME desktop.

Summary

Interactive shell scripts have a reputation for being dull and boring. You can change that by using a few different techniques and tools available on most Linux systems. First, you can create menu systems for your interactive scripts by using the `case` command and shell script functions.

The `case` command allows you to paint a menu, using the standard `echo` command, and read a response from the user, using the `read` command. The `case` command then selects the appropriate shell script function based on the value entered.

The `dialog` program provides several prebuilt text widgets for creating Windows-like objects on a text-based terminal emulator. You can create dialog boxes for displaying text, entering text, and choosing files and dates by using the `dialog` program. This helps bring even more life to your shell script.

If you're running your shell scripts in a graphical X Windows environment, you can utilize even more tools in your interactive scripts. For the KDE desktop, there's the `kdialog` program. This program provides simple commands to create windows widgets for all of the basic windows functions. For the GNOME desktop, there are the `gdialog` and `zenity` programs. Each of these programs provides window widgets that blend into the GNOME desktop just like a real Windows application.

The next chapter dives into the subject of editing and manipulating text data files. Often the biggest use of shell scripts revolves around parsing and displaying data in text files such as log and error files. The Linux environment includes two very useful tools, `sed` and `gawk`, for working with text data in your shell scripts. The next chapter introduces you to these tools, and shows the basics of how to use them.

Chapter 18

Introducing sed and gawk

In This Chapter

- Text manipulation
- The sed editor basics

By far, one of the most common functions that people use shell scripts for is to work with text files. Between examining log files, reading configuration files, and handling data elements, shell scripts can help automate the mundane tasks of manipulating any type of data contained in text files. However, trying to manipulate the contents of text files using just shell script commands can be somewhat awkward. If you perform any type of data manipulation in your shell scripts, you'll want to become familiar with the `sed` and `gawk` tools available in Linux. These tools can greatly simplify any data-handling tasks you need to perform.

Text Manipulation

Chapter 9 showed you how to edit text files using different

editor programs available in the Linux environment. These editors enable you to easily manipulate text contained in a text file by using simple commands or mouse clicks.

There are times, however, when you'll find yourself wanting to manipulate text in a text file on the fly, without having to pull out a full-fledged interactive text editor. In these situations, it would be useful to have a simple command line editor that could easily format, insert, modify, or delete text elements automatically.

The Linux system provides two common tools for doing just that. This section describes the two most popular command line editors used in the Linux world, `sed` and `gawk`.

The sed Editor

The `sed` editor is called a *stream editor*, as opposed to a normal interactive text editor. In an interactive text editor, such as `vim`, you interactively use keyboard commands to insert, delete, or replace text in the data. A stream editor edits a stream of data based on a set of rules you supply ahead of time, before the editor processes the data.

The `sed` editor can manipulate data in a data stream based on commands you either enter into the command line or store in a command text file. It reads one line of data at a time from the input, matches that data with the supplied editor commands, changes data in the stream as specified in the commands, and then outputs the new data to `STDOUT`. After the stream editor matches all of the commands against a line of data, it reads the next line of data and

repeats the process. After the stream editor processes all of the lines of data in the stream, it terminates.

Because the commands are applied sequentially line by line, the `sed` editor has to make only one pass through the data stream to make the edits. This makes the `sed` editor much faster than an interactive editor, and allows you to quickly make changes to data in a file on-the-fly.

The format for using the `sed` command is:

```
sed options script file
```

The options parameters allow you to customize the behavior of the `sed` command, and include the options shown in [Table 18.1](#).

Table 18.1 The `sed` Command Options

Option	Description
<code>-e script</code>	Add commands specified in the script to the commands run while processing the input.
<code>-f file</code>	Add the commands specified in the file to the commands run while processing the input.
<code>-n</code>	Don't produce output for each command, but wait for the <code>print</code> command.

The `script` parameter specifies a single command to apply against the stream data. If more than one command is required, you must use either the `-e` option to specify them in the command line or the `-f` option to specify them in a separate file. Numerous commands are available for manipulating data. We'll examine some of the basic commands used by the `sed` editor later in this chapter, and

then look at some of the more advanced commands in Chapter 20.

Defining an Editor Command in the Command Line

By default, the `sed` editor applies the specified commands to the `STDIN` input stream. This allows you to pipe data directly to the `sed` editor for processing. Here's a quick example demonstrating how to do this:

```
$ echo "This is a test" | sed 's/test/big test/'
This is a big test
$
```

This example uses the `s` command in the `sed` editor. The `s` command substitutes a second text string for the first text string pattern specified between the forward slashes. In this example, the words `big test` were substituted for the word `test`.

When you run this example, it should display the results almost instantaneously. That's the power of using the `sed` editor. You can make multiple edits to data in about the same time it takes for some of the interactive editors just to start up.

Of course, this simple test just edited one line of data. You should get the same speedy results when editing complete files of data:

```
$ cat data1
The quick brown fox jumps over the lazy dog.
The quick brown fox jumps over the lazy dog.
The quick brown fox jumps over the lazy dog.
The quick brown fox jumps over the lazy dog.
```

```
$ sed 's/dog/cat/' datal
The quick brown fox jumps over the lazy cat.
The quick brown fox jumps over the lazy cat.
The quick brown fox jumps over the lazy cat.
The quick brown fox jumps over the lazy cat.
$
```

The `sed` command executes and returns the data almost instantaneously. As it processes each line of data, the results are displayed. You'll start seeing results before the `sed` editor completes processing the entire file.

It's important to note that the `sed` editor doesn't modify the data in the text file itself. It only sends the modified text to `STDOUT`. If you look at the text file, it still contains the original data:

```
$ cat datal
The quick brown fox jumps over the lazy dog.
The quick brown fox jumps over the lazy dog.
The quick brown fox jumps over the lazy dog.
The quick brown fox jumps over the lazy dog.
$
```

Using Multiple Editor Commands in the Command Line

To execute more than one command from the `sed` command line, just use the `-e` option:

```
$ sed -e 's/brown/green/; s/dog/cat/' datal
The quick green fox jumps over the lazy cat.
The quick green fox jumps over the lazy cat.
The quick green fox jumps over the lazy cat.
The quick green fox jumps over the lazy cat.
$
```

Both commands are applied to each line of data in the file. The commands must be separated with a semicolon,

and there shouldn't be any spaces between the end of the command and the semicolon.

Instead of using a semicolon to separate the commands, you can use the secondary prompt in the bash shell. Just enter the first single quotation mark to open the script, and bash will continue to prompt you for more commands until you enter the closing quotation mark:

```
$ sed -e `
> s/brown/green/
> s/fox/elephant/
> s/dog/cat/' data1
The quick green elephant jumps over the lazy cat.
The quick green elephant jumps over the lazy cat.
The quick green elephant jumps over the lazy cat.
The quick green elephant jumps over the lazy cat.
$
```

You must remember to finish the command on the same line that the closing single quotation mark appears. Once the bash shell detects the closing quotation mark, it will process the command. Once it starts, the `sed` command applies each command you specified to each line of data in the text file.

Reading Editor Commands from a File

Finally, if you have lots of `sed` commands you want to process, it is often easier to just store them in a separate file. Use the `-f` option to specify the file in the `sed` command:

```
$ cat script1
s/brown/green/
s/fox/elephant/
s/dog/cat/
$
$ sed -f script1 data1
The quick green elephant jumps over the lazy cat.
```



```
The quick green elephant jumps over the lazy cat.  
The quick green elephant jumps over the lazy cat.  
The quick green elephant jumps over the lazy cat.  
$
```

In this case, you don't put a semicolon after each command. The `sed` editor knows that each line contains a separate command. As with entering commands on the command line, the `sed` editor reads the commands from the specified file and applies them to each line in the data file.

We'll be looking at some other `sed` editor commands that will come in handy for manipulating data in the “The `sed` Editor Basics” section. Before that, let's take a quick look at the other Linux data editor.

The gawk Program

While the `sed` editor is a handy tool for modifying text files on-the-fly, it has its limitations. Often you need a more advanced tool for manipulating data in a file, one that provides a more programming-like environment allowing you to modify and reorganize data in a file. This is where `gawk` comes in.

The `gawk` program is the GNU version of the original `awk` program in Unix. The `gawk` program takes stream editing one step further than the `sed` editor by providing a programming language instead of just editor commands. Within the `gawk` programming language, you can do the following:

- Define variables to store data.
- Use arithmetic and string operators to operate on

data.

- Use structured programming concepts, such as if-then statements and loops, to add logic to your data processing.
- Generate formatted reports by extracting data elements within the data file and repositioning them in another order or format.

The `gawk` program's report-generating capabilities are often used for extracting data elements from large bulky text files and formatting them into a readable report. The perfect example of this is formatting log files. Trying to pore through lines of errors in a log file can be difficult. The `gawk` program allows you to filter just the data elements you want to view from the log file, and then format them in a manner that makes reading the important data easier.

The gawk Command Format

The basic format of the `gawk` program is as follows:

`gawk options program file`

[Table 18.2](#) shows the options available with the `gawk` program.

Table 18.2 The gawk Options

Option	Description
<code>-F fs</code>	Specify a file separator for delineating data fields in a line.
<code>-f file</code>	Specify a file name to read the program from.
<code>-v</code> <code>var=value</code>	Define a variable and default value used in the <code>gawk</code> program.

<code>-mf N</code>	Specify the maximum number of fields to process in the data file.
<code>-mr N</code>	Specify the maximum record size in the data file.
<code>-W keyword</code>	Specify the compatibility mode or warning level for <code>gawk</code> .

The command line options provide an easy way to customize features in the `gawk` program. We'll be looking more closely at these as we explore `gawk`.

The power of `gawk` is in the program script. You can write scripts to read the data within a text line and then manipulate and display the data to create any type of output report.

Reading the Program Script from the Command Line

A `gawk` program script is defined by opening and closing braces. You must place script commands between the two braces. Because the `gawk` command line assumes that the script is a single text string, you must also enclose your script in single quotation marks. The following is an example of a simple `gawk` program script specified on the command line:

```
$ gawk '{print "Hello John!"}'
```

The program script defines a single command, the `print` command. The `print` command does what it says: it prints text to `STDOUT`. If you try running this command, you'll be somewhat disappointed, as nothing will happen right away. Because no file name was defined in the command line,

the `gawk` program retrieves data from `STDIN`. When you run the program, it just waits for text to come in via `STDIN`.

If you type a line of text and press the Enter key, `gawk` runs the text through the program script:

```
$ gawk '{print "Hello John!"}'
This is a test
Hello John!
hello
Hello John!
This is another test
Hello John!

$
```

Just like the `sed` editor, the `gawk` program executes the program script on each line of text available in the data stream. Because the program script is set to display a fixed text string, no matter what text you enter in the data stream, you'll get the same text output.

To terminate the `gawk` program, you must signal that the data stream has ended. The bash shell provides a key combination to generate an End-of-File (EOF) character. The `Ctrl+D` key combination generates an EOF character in bash. Using that key combination terminates the `gawk` program and returns you to a command line interface prompt.

Using Data Field Variables

One of the primary features of `gawk` is its ability to manipulate data in the text file. It does this by automatically assigning a variable to each data element in a line. By default, `gawk` assigns the following variables to each data

field it detects in the line of text:

- `$0` represents the entire line of text.
- `$1` represents the first data field in the line of text.
- `$2` represents the second data field in the line of text.
- `$n` represents the *n*th data field in the line of text.

Each data field is determined in a text line by a *field separation character*. When `gawk` reads a line of text, it delineates each data field using the defined field separation character. The default field separation character in `gawk` is any whitespace character (such as the tab or space characters).

The following is an example `gawk` program that reads a text file and displays only the first data field value:

```
$ cat data3
One line of test text.
Two lines of test text.
Three lines of test text.
$
$ gawk '{print $1}' data3
One
Two
Three
$
```

This program uses the `$1` field variable to display only the first data field for each line of text.

If you're reading a file that uses a different field separation character, you can specify it by using the `-F` option:

```
$ gawk -F: '{print $1}' /etc/passwd
root
daemon
bin
sys
sync
```

```
games
man
lp
mail
...
```

This short program displays the first data field in the password file on the system. Because the `/etc/passwd` file uses a colon to separate the data fields, if you want to separate out each data element you must specify it as the field separation character in the `gawk` options.

Using Multiple Commands in the Program Script

A programming language wouldn't be very useful if you could only execute one command. The `gawk` programming language allows you to combine commands into a normal program. To use multiple commands in the program script specified on the command line, just place a semicolon between each command:

```
$ echo "My name is Rich" | gawk '{ $4="Christine"; print $0 }'
My name is Christine
$
```

The first command assigns a value to the `$4` field variable. The second command then prints the entire data field. Notice from the output that the `gawk` program replaced the fourth data field in the original text with the new value.

You can also use the secondary prompt to enter your program script commands one line at a time:

```
$ gawk '{
> $4="testing"
> print $0 }'
This is not a good test.
```

```
This is not testing good test.
```

```
$
```

After you open the single quotation mark, the bash shell provides the secondary prompt to prompt you for more data. You can add your commands one at a time on each line until you enter the closing single quotation mark. To exit the program, just press the Ctrl+D key combination to signal the end of the data.

Reading the Program from a File

As with the `sed` editor, the `gawk` editor allows you to store your programs in a file and refer to them in the command line:

```
$ cat script2
{ print $1 "'s home directory is " $6 }
$
$ gawk -F: -f script2 /etc/passwd
root's home directory is /root
daemon's home directory is /usr/sbin
...
Samantha's home directory is /home/Samantha
Timothy's home directory is /home/Timothy
Christine's home directory is /home/Christine
$
```

The `script2` program script uses the `print` command again to print the `/etc/passwd` file's home directory data field (field variable `$6`), and the `userid` data field (field variable `$1`).

You can specify multiple commands in the program file. To do so, just place each command on a separate line. There's no need to use semicolons:

```
$
$ cat script3
{
text = "'s home directory is "
print $1 text $6
}
```

```
$  
$ gawk -F: -f script3 /etc/passwd  
root's home directory is /root  
daemon's home directory is /usr/sbin  
...  
Samantha's home directory is /home/Samantha  
Timothy's home directory is /home/Timothy  
Christine's home directory is /home/Christine  
$
```

The `script3` program script defines a variable to hold a text string used in the `print` command. You'll notice that `gawk` programs don't use a dollar sign when referencing a variable's value, as a shell script does.

Running Scripts Before Processing Data

The `gawk` program also allows you to specify when the program script is run. By default, `gawk` reads a line of text from the input and then executes the program script on the data in the line of text. Sometimes you may need to run a script before processing data, such as to create a header section for a report. The `BEGIN` keyword is used to accomplish this. It forces `gawk` to execute the program script specified after the `BEGIN` keyword, before `gawk` reads the data:

```
$ gawk 'BEGIN {print "Hello World!"}'  
Hello World!  
$
```

This time the `print` command displays the text before reading any data. However, after it displays the text, it quickly exits, without waiting for any data.

The reason for this is that the `BEGIN` keyword only applies the specified script before it processes any data. If you want to process data with a normal program script, you must define the program using another script section:

```
$
$ cat data4
Line 1
Line 2
Line 3
$
$ gawk 'BEGIN { print "The data4 File Contents:" } { print $0 }' data4
The data4 File Contents:
Line 1
Line 2
Line 3
$
```

Now after `gawk` executes the `BEGIN` script, it uses the second script to process any file data. Be careful when doing this—notice that both of the scripts are still considered one text string on the `gawk` command line. You need to place your single quotation marks accordingly.

Running Scripts After Processing Data

Like the `BEGIN` keyword, the `END` keyword allows you to specify a program script that `gawk` executes after reading the data:

```
$
$ gawk 'BEGIN { print "The data4 File Contents:" } { print $0 }
  END { print "End of File" }' data4
The data4 File Contents:
Line 1
Line 2
Line 3
End of File
$
```

When the `gawk` program is done printing out the file

contents, it executes the commands in the END script. This is a great technique to use to add footer data to reports after all the normal data has been processed.

You can put all of these elements together into a nice little program script file to create a full report from a simple data file:

```
$ cat script4
BEGIN {
print "The latest list of users and shells"
print "  Userid      Shell"
print "-----      -"
FS=":"
}

{
print $1 "      " $7"
}

END {
print "This concludes the listing"
}
$
```

This script uses the BEGIN script to create a header section for the report. It also defines a special variable called FS. This is yet another way to define the field separation character. This way you don't have to depend on the script's user to define the field separation character in the command line options.

Here's a somewhat truncated output from running this gawk program script:

```
$ gawk -f script4 /etc/passwd
The latest list of users and shells
  Userid      Shell
-----      -
root         /bin/bash
daemon       /bin/sh
bin          /bin/sh
...
```

```
Samantha      /bin/bash
Timothy       /bin/sh
Christine     /bin/sh
This concludes the listing
$
```

As expected, the `BEGIN` script created the header text, the program script processed the information from the specified data file (the `/etc/passwd` file), and the `END` script produced the footer text.

This gives you a small taste of the power available when you use simple `gawk` scripts. Chapter 21 describes some more basic programming principles available for your `gawk` scripts, along with some even more advanced programming concepts you can use in your `gawk` program scripts to create professional looking reports from even the most cryptic data files.

The sed Editor Basics

The key to successfully using the `sed` editor is to know its myriad of commands and formats, which help you customize your text editing. This section describes some of the basic commands and features you can incorporate into your script to start using the `sed` editor.

More Substitution Options

You've already seen how to use the `s` command to substitute new text for the text in a line. However, a few additional options are available for the `substitute` command

that can help make your life easier.

Substitution Flags

There's a caveat to how the `substitute` command replaces matching patterns in the text string. Watch what happens in this example:

```
$ cat data5
This is a test of the test script.
This is the second test of the test script.
$
$ sed 's/test/trial/' data5
This is a trial of the test script.
This is the second trial of the test script.
$
```

The `substitute` command works fine in replacing text in multiple lines, but by default, it only replaces the first occurrence in each line. To get the `substitute` command to work on different occurrences of the text, you must use a *substitution flag*. The substitution flag is set after the substitution command strings:

`s/pattern/replacement/flags`

There are four types of substitution flags available:

- A number, indicating the pattern occurrence for which new text should be substituted.
- `g`—Indicates that new text should be substituted for all occurrences of the existing text.
- `p`—Indicates that the contents of the original line should be printed.
- `w file`—Write the results of the substitution to a file.

In the first type of substitution, you can specify which

occurrence of the matching pattern the `sed` editor should substitute new text for:

```
$ sed 's/test/trial/2' data5
This is a test of the trial script.
This is the second test of the trial script.
$
```

As a result of specifying a `2` as the substitution flag, the `sed` editor only replaces the pattern in the second occurrence in each line. The `g` substitution flag enables you to replace every occurrence of the pattern in the text:

```
$ sed 's/test/trial/g' data5
This is a trial of the trial script.
This is the second trial of the trial script.
$
```

The `p` substitution flag prints a line that contains a matching pattern in the substitute command. This is most often used in conjunction with the `-n sed` option:

```
$ cat data6
This is a test line.
This is a different line.
$
$ sed -n 's/test/trial/p' data5
This is a trial line.
$
```

The `-n` option suppresses output from the `sed` editor. However, the `p` substitution flag outputs any line that has been modified. Using the two in combination produces output only for lines that have been modified by the substitute command.

The `w` substitution flag produces the same output but stores the output in the specified file:

```
$ sed 's/test/trial/w test' data6
This is a trial line.
This is a different line.
```

```
$ cat test
This is a trial line.
$
```

The normal output of the `sed` editor appears in `STDOUT`, but only the lines that include the matching pattern are stored in the specified output file.

Replacement Characters

There are times when you run across characters in text strings that aren't easy to use in the substitution pattern. One popular example in the Linux world is the forward slash.

Substituting pathnames in a file can get awkward. For example, if you wanted to substitute the C shell for the bash shell in the `/etc/passwd` file, you'd have to do this:

```
$ sed 's/\\bin\\bash/\\bin\\csh/' /etc/passwd
```

Because the forward slash is used as the string delimiter, you must use a backslash to escape it if it appears in the pattern text. This often leads to confusion and mistakes.

To solve this problem, the `sed` editor allows you to select a different character for the string delimiter in the substitute command:

```
$ sed 's!/bin/bash!/bin/csh!' /etc/passwd
```

In this example, the exclamation point is used for the string delimiter, making the pathnames much easier to read and understand.

Using Addresses

By default, the commands you use in the `sed` editor apply to all lines of the text data. If you only want to apply a command to a specific line, or a group of lines, you must use *line addressing*.

There are two forms of line addressing in the `sed` editor:

- A numeric range of lines
- A text pattern that filters out a line

Both forms use the same format for specifying the address:

```
[address]command
```

You can also group more than one command together for a specific address:

```
address {  
    command1  
    command2  
    command3  
}
```

The `sed` editor applies each of the commands you specify only to lines that match the address specified.

This section demonstrates using both of these addressing techniques in your `sed` editor scripts.

Numeric Line Addressing

When using numeric line addressing, you reference lines using their line position in the text stream. The `sed` editor assigns the first line in the text stream as line number one and continues sequentially for each new line.

The address you specify in the command can be a single line number or a range of lines specified by a starting line number, a comma, and an ending line number. Here's an example of specifying a line number to which the `sed` command will be applied:

```
$ sed '2s/dog/cat/' data1
The quick brown fox jumps over the lazy dog
The quick brown fox jumps over the lazy cat
The quick brown fox jumps over the lazy dog
The quick brown fox jumps over the lazy dog
$
```

The `sed` editor modified the text only in line two per the address specified. Here's another example, this time using a range of line addresses:

```
$ sed '2,3s/dog/cat/' data1
The quick brown fox jumps over the lazy dog
The quick brown fox jumps over the lazy cat
The quick brown fox jumps over the lazy cat
The quick brown fox jumps over the lazy dog
$
```

If you want to apply a command to a group of lines starting at some point within the text, but continuing to the end of the text, you can use the special address, the dollar sign:

```
$ sed '2,$s/dog/cat/' data1
The quick brown fox jumps over the lazy dog
The quick brown fox jumps over the lazy cat
The quick brown fox jumps over the lazy cat
The quick brown fox jumps over the lazy cat
$
```

Because you may not know how many lines of data are in the text, the dollar sign often comes in handy.

Using Text Pattern Filters

The other method of restricting which lines a command applies to is a bit more complicated. The `sed` editor allows you to specify a text pattern that it uses to filter lines for the command. The format for this is:

```
/pattern/command
```

You must encapsulate the *pattern* you specify in forward slashes. The `sed` editor applies the command only to lines that contain the text pattern that you specify.

For example, if you want to change the default shell for only the user Samantha, you'd use the `sed` command:

```
$
$ grep Samantha /etc/passwd
Samantha:x:1001:1002:Samantha,4,,:/home/Samantha:/bin/bash
$
$ sed '/Samantha/s/bash/csh/' /etc/passwd
root:x:0:0:root:/root:/bin/bash
...
Samantha:x:1001:1002:Samantha,4,,:/home/Samantha:/bin/csh
Timothy:x:1002:1005:/home/Timothy:/bin/sh
Christine:x:1003:1006:/home/Christine:/bin/sh
$
```

The command was only applied to the line with the matching text pattern. While using a fixed text pattern may be useful for filtering specific values, as in the `userid` example, it's somewhat limited in what you can do with it. The `sed` editor uses a feature called *regular expressions* in text patterns to allow you to create patterns that get pretty involved.

Regular expressions allow you to create advanced text pattern-matching formulas to match all sorts of data. These formulas combine a series of wildcard characters, special characters, and fixed text characters to produce a concise

pattern that can match just about any text situation. Regular expressions are one of the scarier parts of shell script programming and Chapter 19 covers them in great detail.

Grouping Commands

If you need to perform more than one command on an individual line, group the commands together using braces. The `sed` editor will process each command listed on the address line(s):

```
$ sed '2{
> s/fox/elephant/
> s/dog/cat/
> }' data1
The quick brown fox jumps over the lazy dog.
The quick brown elephant jumps over the lazy cat.
The quick brown fox jumps over the lazy dog.
The quick brown fox jumps over the lazy dog.
$
```

Both commands are processed against the address. And of course, you can also specify an address range before the grouped commands:

```
$ sed '3,${
> s/brown/green/
> s/lazy/active/
> }' data1
The quick brown fox jumps over the lazy dog.
The quick brown fox jumps over the lazy dog.
The quick green fox jumps over the active dog.
The quick green fox jumps over the active dog.
$
```

The `sed` editor applies all of the commands to all of the lines in the address range.

Deleting Lines

The text substitution command isn't the only command available in the `sed` editor. If you need to delete specific lines of text in a text stream, there's the `delete` command.

The `delete` command, `d`, pretty much does what it says. It will delete any text lines that match the addressing scheme supplied. Be careful with the `delete` command, because if you forget to include an addressing scheme, all of the lines will be deleted from the stream:

```
$ cat data1
The quick brown fox jumps over the lazy dog
The quick brown fox jumps over the lazy dog
The quick brown fox jumps over the lazy dog
The quick brown fox jumps over the lazy dog
$ sed 'd' data1
$
```

The `delete` command is obviously most useful when used in conjunction with a specified address. This allows you to delete specific lines of text from the data stream, either by line number:

```
$ sed '3d' data7
This is line number 1.
This is line number 2.
This is line number 4.
$
```

or by a specific range of lines:

```
$ sed '2,3d' data7
This is line number 1.
This is line number 4.
$
```

or by using the special end-of-file character:

```
$ sed '3,$d' data7
This is line number 1.
This is line number 2.
$
```

The pattern-matching feature of the `sed` editor also applies to the `delete` command:

```
$ sed '/number 1/d' data7
This is line number 2.
This is line number 3.
This is line number 4.
$
```

The `sed` editor removes the line containing text that matches the pattern you specify.

Note

Remember that the `sed` editor doesn't touch the original file. Any lines you delete are only gone from the output of the `sed` editor. The original file still contains the “deleted” lines.

You can also delete a range of lines using two text patterns, but be careful if you do this. The first pattern you specify “turns on” the line deletion, and the second pattern “turns off” the line deletion. The `sed` editor deletes any lines between the two specified lines (including the specified lines):

```
$ sed '/1/,/3/d' data6
This is line number 4.
$
```

In addition, you need to be careful, as the delete feature will “turn on” whenever the `sed` editor detects the start pattern in the data stream. This may produce an unexpected result:

```
$ cat data8
This is line number 1.
This is line number 2.
This is line number 3.
```

```
This is line number 1.  
This is line number 1 again.  
This is text you want to keep.  
This is the last line in the file.  
$  
$ sed '/1/,/3/d' data7  
This is line number 4.  
$
```

The second occurrence of a line with the number 1 in it triggered the `delete` command again, deleting the rest of the lines in the data stream, as the stop pattern wasn't recognized. Of course, the other obvious problem occurs if you specify a stop pattern that never appears in the text:

```
$ sed '/1/,/5/d' data8  
$
```

Because the `delete` features “turned on” at the first pattern match, but never found the end pattern match, the entire data stream was deleted.

Inserting and Appending Text

As you would expect, like any other editor, the `sed` editor allows you to insert and append text lines to the data stream. The difference between the two actions can be confusing:

- The `insert` command (`i`) adds a new line before the specified line.
- The `append` command (`a`) adds a new line after the specified line.

What's confusing about these two commands is their formats. You can't use these commands on a single command line. You must specify the line to insert or append

on a separate line by itself. The format for doing this is:

```
sed '[address]command\  
new line'
```

The text in *new line* appears in the `sed` editor output in the place you specify. Remember that when you use the `insert` command, the text appears before the data stream text:

```
$ echo "Test Line 2" | sed 'i\Test Line 1'  
Test Line 1  
Test Line 2  
$
```

And when you use the `append` command, the text appears after the data stream text:

```
$ echo "Test Line 2" | sed 'a\Test Line 1'  
Test Line 2  
Test Line 1  
$
```

When you use the `sed` editor from the command line interface prompt, you'll get the secondary prompt to enter the new line data. You must complete the `sed` editor command on this line. Once you enter the ending single quotation mark, the bash shell will process the command:

```
$ echo "Test Line 2" | sed 'i\  
> Test Line 1'  
Test Line 1  
Test Line 2  
$
```

This works well for adding text before or after the text in the data stream, but what about adding text inside the data stream?

To insert or append data inside the data stream lines, you must use addressing to tell the `sed` editor where you want the data to appear. You can specify only a single line

address when using these commands. You can match either a numeric line number or a text pattern, but you cannot use a range of addresses. This is logical, because you can only insert or append before or after a single line, and not a range of lines.

The following is an example of inserting a new line before line 3 in the data stream:

```
$ sed '3i\  
> This is an inserted line.' data7  
This is line number 1.  
This is line number 2.  
This is an inserted line.  
This is line number 3.  
This is line number 4.  
$
```

The following is an example of appending a new line after line 3 in the data stream:

```
$ sed '3a\  
>This is an appended line.' data7  
This is line number 1.  
This is line number 2.  
This is line number 3.  
This is an appended line.  
This is line number 4.  
$
```

This uses the same process as the `insert` command; it just places the new text line after the specified line number. If you have a multiline data stream, and you want to append a new line of text to the end of a data stream, just use the dollar sign, which represents the last line of data:

```
$ sed '$a\  
> This is a new line of text.' data7  
This is line number 1.  
This is line number 2.  
This is line number 3.  
This is line number 4.  
This is a new line of text.
```

The same idea applies if you want to add a new line at the beginning of the data stream. Just insert a new line before line number one.

To insert or append more than one line of text, you must use a backslash on each line of new text until you reach the last text line where you want to insert or append text:

```
$ sed '1i\  
> This is one line of new text.\  
> This is another line of new text.' data7  
This is one line of new text.  
This is another line of new text.  
This is line number 1.  
This is line number 2.  
This is line number 3.  
This is line number 4.  
$
```

Both of the specified lines are added to the data stream.

Changing Lines

The `change` command allows you to change the contents of an entire line of text in the data stream. It works the same way as the `insert` and `append` commands, in that you must specify the new line separately from the rest of the `sed` command:

```
$ sed '3c\  
> This is a changed line of text.' data7  
This is line number 1.  
This is line number 2.  
This is a changed line of text.  
This is line number 4.  
$
```

In this example, the `sed` editor changes the text in line number 3. You can also use a text pattern for the address:


```
$ sed '/number 3/c\  
> This is a changed line of text.' data7  
This is line number 1.  
This is line number 2.  
This is a changed line of text.  
This is line number 4.  
$
```

The text pattern `change` command will change any line of text in the data stream that it matches.

```
$ sed '/number 1/c\  
> This is a changed line of text.' data8  
This is a changed line of text.  
This is line number 2.  
This is line number 3.  
This is line number 4.  
This is a changed line of text.  
This is yet another line.  
This is the last line in the file.  
$
```

You can use an address range in the `change` command, but the results may not be what you expect:

```
$ sed '2,3c\  
> This is a new line of text.' data7  
This is line number 1.  
This is a new line of text.  
This is line number 4.  
$
```

Instead of changing both lines with the text, the `sed` editor uses the single line of text to replace both lines.

The transform Command

The `transform` command (`y`) is the only `sed` editor command that operates on a single character. The `transform` command uses the format:

```
[address]y/inchars/outchars/
```

The `transform` command performs a one-to-one mapping of the `inchars` and the `outchars` values. The first character in `inchars` is converted to the first character in `outchars`. The second character in `inchars` is converted to the second character in `outchars`. This mapping continues throughout the length of the specified characters. If the `inchars` and `outchars` are not the same length, the `sed` editor will produce an error message.

A simple example of using the `transform` command is:

```
$ sed 'y/l23/789/' data8
This is line number 7.
This is line number 8.
This is line number 9.
This is line number 4.
This is line number 7 again.
This is yet another line.
This is the last line in the file.
$
```

As you can see from the output, each instance of the characters specified in the `inchars` pattern has been replaced by the character in the same position in the `outchars` pattern.

The `transform` command is a global command; that is, it performs the transformation on any character found in the text line automatically, without regard to the occurrence:

```
$ echo "This 1 is a test of 1 try." | sed 'y/123/456/'
This 4 is a test of 4 try.
$
```

The `sed` editor transformed both instances of the matching character 1 in the text line. You can't limit the transformation to a specific occurrence of the character.

Printing Revisited

The “More Substitution Options” section showed you how to use the `p` flag with the substitution command to display lines that the `sed` editor changed. There are three commands that also can be used to print information from the data stream:

- The lowercase `p` command to print a text line
- The equal sign (`=`) command to print line numbers
- The `l` (lowercase L) command to list a line

The following sections look at each of these three printing commands in the `sed` editor.

Printing Lines

Like the `p` flag in the substitution command, the `p` command prints a line in the `sed` editor output. On its own, there's not much excitement:

```
$ echo "this is a test" | sed 'p'
this is a test
this is a test
$
```

All it does is print the data text that you already know is there. The most common use for the `print` command is printing lines that contain matching text from a text pattern:

```
$ sed -n '/number 3/p' data7
This is line number 3.
$
```

By using the `-n` option on the command line, you can suppress all of the other lines and only print the line that contains the matching text pattern.

You can also use this as a quick way to print a subset of

lines in a data stream:

```
$ sed -n '2,3p' data7
This is line number 2.
This is line number 3.
$
```

You can also use the `print` command when you need to see a line before it gets altered, such as with the `substitution` or `change` command. You can create a script that displays the line before it's changed:

```
$ sed -n '/3/{
p
s/line/test/p
}' data7
This is line number 3.
This is test number 3.
$
```

This `sed` editor command searches for lines that contain the number 3, and then executes two commands. First, the script uses the `p` command to print the original version of the line; then it uses the `s` command to substitute text, along with the `p` flag to print the resulting text. The output shows both the original line text and the new line text.

Printing Line Numbers

The `equal sign` command prints the current line number for the line within the data stream. Line numbers are determined by using the newline character in the data stream. Each time a newline character appears in the data stream, the `sed` editor assumes that it terminates a line of text:

```
$ sed '=' data1
1
```

```
The quick brown fox jumps over the lazy dog.  
2  
The quick brown fox jumps over the lazy dog.  
3  
The quick brown fox jumps over the lazy dog.  
4  
The quick brown fox jumps over the lazy dog.  
$
```

The `sed` editor prints the line number before the actual line of text. The `equal` sign command comes in handy if you're searching for a specific text pattern in the data stream:

```
$ sed -n '/number 4/{  
=  
P  
}' data7  
4  
This is line number 4.  
$
```

By using the `-n` option you can have the `sed` editor display only both the line number and text for the line that contains the matching text pattern.

Listing Lines

The `list` command (`l`) allows you to print both the text and nonprintable ASCII characters in a data stream. Any nonprintable characters are shown using either their octal values, preceded by a backslash, or the standard C-style nomenclature for common nonprintable characters, such as `\t` for tab characters:

```
$ cat data9  
This line contains tabs.  
$  
$ sed -n 'l' data9  
This\tline\tcontains\ttabs.$  
$
```

The tab character locations are shown with the `\t` nomenclature. The dollar sign at the end of the line indicates the newline character. If you have a data stream that contains an escape character, the `list` command displays it using the octal code:

```
$ cat data10
This line contains an escape character
$
$ sed -n 'l' data10
This line contains an escape character \033[44m$
$
```

The `data10` text file contains an escape control code (see Chapter 17) to change the display color. When you use the `cat` command to display the text file, you don't see the escape control code; it just changes the color of the display.

However, using the `list` command, you can display the escape control code used. The `\033` is the octal value of the ASCII code for the escape key.

Using Files with sed

The `substitution` command contains flags that allow you to work with files. There are also regular `sed` editor commands that let you do that without having to substitute text.

Writing to a File

The `w` command is used to write lines to a file. The format for the `w` command is:

```
[address]w filename
```

The `filename` can be specified as either a relative or absolute pathname, but in either case, the person running the `sed` editor must have write permissions for the file. The address can be any type of addressing method used in `sed`, such as a single line number, a text pattern, or a range of line numbers or text patterns.

The following is an example that prints only the first two lines of a data stream to a text file:

```
$ sed '1,2w test' data7
This is line number 1.
This is line number 2.
This is line number 3.
This is line number 4.
$
$ cat test
This is line number 1.
This is line number 2.
$
```

Of course, if you don't want the lines to display on `STDOUT`, you can use the `-n` option for the `sed` command.

This is a great tool to use if you need to create a data file from a master file on the basis of common text values, such as those in a mailing list:

```
$ cat datall
Blum, Katie           Chicago, IL
Mullen, Riley         West Lafayette, IN
Snell, Haley          Ft. Wayne, IN
Woenker, Matthew      Springfield, IL
Wisecarver, Emma      Grant Park, IL
$
$ sed -n '/IN/w INcustomers' datall
$
$ cat INcustomers
Mullen, Riley         West Lafayette, IN
Snell, Haley          Ft. Wayne, IN
$
```

The `sed` editor writes to a destination file only the data

lines that contain the text pattern.

Reading Data from a File

You've already seen how to insert data into and append text to a data stream from the `sed` command line. The `read` command (`r`) allows you to insert data contained in a separate file.

The format of the `read` command is:

```
[address]r filename
```

The `filename` parameter specifies either an absolute or relative pathname for the file that contains the data. You can't use a range of addresses for the `read` command. You can only specify a single line number or text pattern address. The `sed` editor inserts the text from the file after the address.

```
$ cat data12
This is an added line.
This is the second added line.
$
$ sed '3r data12' data7
This is line number 1.
This is line number 2.
This is line number 3.
This is an added line.
This is the second added line.
This is line number 4.
$
```

The `sed` editor inserts into the data stream all of the text lines in the data file. The same technique works when using a text pattern address:

```
$ sed '/number 2/r data12' data7
This is line number 1.
This is line number 2.
This is an added line.
```



```
This is the second added line.  
This is line number 3.  
This is line number 4.  
$
```

If you want to add text to the end of a data stream, just use the dollar sign address symbol:

```
$ sed '$r data12' data7  
This is line number 1.  
This is line number 2.  
This is line number 3.  
This is line number 4.  
This is an added line.  
This is the second added line.  
$
```

A cool application of the `read` command is to use it in conjunction with a `delete` command to replace a placeholder in a file with data from another file. For example, suppose that you had a form letter stored in a text file that looked like this:

```
$ cat letter  
Would the following people:  
LIST  
please report to the office.  
$
```

The form letter uses the generic placeholder `LIST` in place of a list of people. To insert the list of people after the placeholder, all you need to do is use the `read` command. However, this still leaves the placeholder text in the output. To remove that, just use the `delete` command. The result looks like this:

```
$ sed '/LIST/{  
> r data11  
> d  
> }' letter  
Would the following people:  
Blum, Katie           Chicago, IL  
Mullen, Riley         West Lafayette, IN
```

```
Snell, Haley      Ft. Wayne, IN
Woenker, Matthew  Springfield, IL
Wisecarver, Emma  Grant Park, IL
please report to the office.
$
```

Now the placeholder text is replaced with the list of names from the data file.

Summary

While shell scripts can do a lot of work on their own, it's often difficult to manipulate data with just a shell script. Linux provides two handy utilities to help out with handling text data. The `sed` editor is a stream editor that quickly processes data on the fly as it reads it. You must provide the `sed` editor with a list of editing commands, which it applies to the data.

The `gawk` program is a utility from the GNU organization that mimics and expands on the functionality of the Unix `awk` program. The `gawk` program contains a built-in programming language that you can use to write scripts to handle and process data. You can use the `gawk` program to extract data elements from large data files and output them in just about any format you desire. This makes processing large log files a snap, as well as creating custom reports from data files.

A crucial element of using both the `sed` and `gawk` programs is knowing how to use regular expressions. Regular expressions are key to creating customized filters for

extracting and manipulating data in text files. The next chapter dives into the often misunderstood world of regular expressions, showing you how to build regular expressions for manipulating all types of data.

Chapter 19

Regular Expressions

In This Chapter

- Defining regular expressions
- Looking at the basics
- Extending our patterns
- Creating expressions

The key to successfully working with the sed editor and the gawk program in your shell script is your comfort using regular expressions. This is not always an easy thing to do, as trying to filter specific data from a large batch of data can (and often does) get complicated. This chapter describes how to create regular expressions in both the sed editor and the gawk program that can filter out just the data you need.

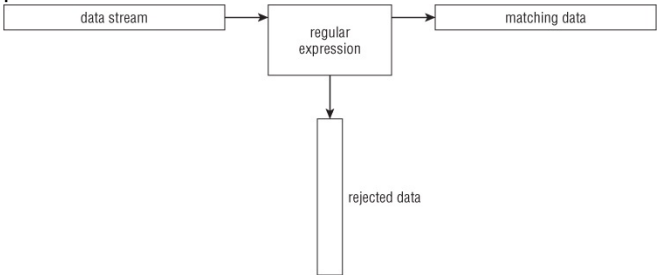
What Are Regular Expressions?

The first step to understanding regular expressions is to define just exactly what they are. This section explains what a regular expression is and describes how Linux uses regular expressions.

A Definition

A regular expression is a pattern template you define that a Linux utility uses to filter text. A Linux utility (such as the sed editor or the gawk program) matches the regular expression pattern against data as that data flows into the utility. If the data matches the pattern, it's accepted for processing. If the data doesn't match the pattern, it's rejected. This is illustrated in [Figure 19.1](#).

Figure 19.1 Matching data against a regular expression pattern



The regular expression pattern makes use of wildcard characters to represent one or more characters in the data

stream. There are plenty of instances in Linux where you can specify a wildcard character to represent data that you don't know about. You've already seen an example of using wildcard characters with the Linux `ls` command for listing files and directories (see Chapter 3).

The asterisk wildcard character allows you to list only files that match a certain criteria. For example:

```
$ ls -al da*
-rw-r--r-- 1 rich    rich          45 Nov 26 12:42 data
-rw-r--r-- 1 rich    rich          25 Dec  4 12:40 data.tst
-rw-r--r-- 1 rich    rich        180 Nov 26 12:42 data1
-rw-r--r-- 1 rich    rich          45 Nov 26 12:44 data2
-rw-r--r-- 1 rich    rich          73 Nov 27 12:31 data3
-rw-r--r-- 1 rich    rich          79 Nov 28 14:01 data4
-rw-r--r-- 1 rich    rich        187 Dec  4 09:45 datatest
$
```

The `da*` parameter instructs the `ls` command to list only the files whose name starts with `da`. There can be any number of characters after the `da` in the filename (including none). The `ls` command reads the information regarding all of the files in the directory but displays only the ones that match the wildcard character.

Regular expression wildcard patterns work in a similar way. The regular expression pattern contains text and/or special characters that define a template for the `sed` editor and the `gawk` program to follow when matching data. There are different special characters you can use in a regular expression to define a specific pattern for filtering data.

Types of Regular Expressions

The biggest problem with using regular expressions is that

there isn't just one set of them. Several different applications use different types of regular expressions in the Linux environment. These include such diverse applications as programming languages (Java, Perl, and Python), Linux utilities (such as the sed editor, the gawk program, and the grep utility), and mainstream applications (such as the MySQL and PostgreSQL database servers).

A regular expression is implemented using a *regular expression engine*. A regular expression engine is the underlying software that interprets regular expression patterns and uses those patterns to match text.

In the Linux world, there are two popular regular expression engines:

- The POSIX Basic Regular Expression (BRE) engine
- The POSIX Extended Regular Expression (ERE) engine

Most Linux utilities at a minimum conform to the POSIX BRE engine specifications, recognizing all of the pattern symbols it defines. Unfortunately, some utilities (such as the sed editor) only conform to a subset of the BRE engine specifications. This is due to speed constraints, as the sed editor attempts to process text in the data stream as quickly as possible.

The POSIX ERE engine is often found in programming languages that rely on regular expressions for text filtering. It provides advanced pattern symbols as well as special symbols for common patterns, such as matching digits, words, and alphanumeric characters. The gawk program

uses the ERE engine to process its regular expression patterns.

Because there are so many different ways to implement regular expressions, it's hard to present a single, concise description of all the possible regular expressions. The following sections discuss the most commonly found regular expressions and demonstrate how to use them in the sed editor and gawk program.

Defining BRE Patterns

The most basic BRE pattern is matching text characters in a data stream. This section demonstrates how you can define text in the regular expression pattern and what to expect from the results.

Plain Text

Chapter 18 demonstrated how to use standard text strings in the sed editor and the gawk program to filter data. Here's an example to refresh your memory:

```
$ echo "This is a test" | sed -n '/test/p'
This is a test
$ echo "This is a test" | sed -n '/trial/p'
$
$ echo "This is a test" | gawk '/test/{print $0}'
This is a test
$ echo "This is a test" | gawk '/trial/{print $0}'
$
```

The first pattern defines a single word, *test*. The sed editor and gawk program scripts each use their own

version of the `print` command to print any lines that match the regular expression pattern. Because the `echo` statement contains the word “test” in the text string, the data stream text matches the defined regular expression pattern, and the `sed` editor displays the line.

The second pattern again defines just a single word, this time the word “trial.” Because the `echo` statement text string doesn't contain that word, the regular expression pattern doesn't match, so neither the `sed` editor nor the `gawk` program prints the line.

You probably already noticed that the regular expression doesn't care where in the data stream the pattern occurs. It also doesn't matter how many times the pattern occurs. Once the regular expression can match the pattern anywhere in the text string, it passes the string along to the Linux utility that's using it.

The key is matching the regular expression pattern to the data stream text. It's important to remember that regular expressions are extremely picky about matching patterns. The first rule to remember is that regular expression patterns are case sensitive. This means they'll only match patterns with the proper case of characters:

```
$ echo "This is a test" | sed -n '/this/p'
$
$ echo "This is a test" | sed -n '/This/p'
This is a test
$
```

The first attempt failed to match because the word “this” doesn't appear in all lowercase in the text string, while the second attempt, which uses the uppercase letter in the

pattern, worked just fine.

You don't have to limit yourself to whole words in the regular expression. If the defined text appears anywhere in the data stream, the regular expression will match the following:

```
$ echo "The books are expensive" | sed -n '/book/p'
The books are expensive
$
```

Even though the text in the data stream is `books`, the data in the stream contains the regular expression `book`, so the regular expression pattern matches the data. Of course, if you try the opposite, the regular expression will fail:

```
$ echo "The book is expensive" | sed -n '/books/p'
$
```

The complete regular expression text didn't appear in the data stream, so the match failed and the sed editor didn't display the text.

You also don't have to limit yourself to single text words in the regular expression. You can include spaces and numbers in your text string as well:

```
$ echo "This is line number 1" | sed -n '/ber 1/p'
This is line number 1
$
```

Spaces are treated just like any other character in the regular expression:

```
$ echo "This is line number1" | sed -n '/ber 1/p'
$
```

If you define a space in the regular expression, it must appear in the data stream. You can even create a regular expression pattern that matches multiple contiguous

spaces:

```
$ cat data1
This is a normal line of text.
This is  a line with too many spaces.
$ sed -n '/ /p' data1
This is  a line with too many spaces.
$
```

The line with two spaces between words matches the regular expression pattern. This is a great way to catch spacing problems in text files!

Special Characters

As you use text strings in your regular expression patterns, there's something you need to be aware of. There are a few exceptions when defining text characters in a regular expression. Regular expression patterns assign a special meaning to a few characters. If you try to use these characters in your text pattern, you won't get the results you were expecting.

The special characters recognized by regular expressions are:

```
.*[]^${}\\+?|()
```

As the chapter progresses, you'll find out just what these special characters do in a regular expression. For now, however, just remember that you can't use these characters by themselves in your text pattern.

If you want to use one of the special characters as a text character, you need to *escape* it. When you escape the special characters, you add a special character in front of it

to indicate to the regular expression engine that it should interpret the next character as a normal text character. The special character that does this is the backslash character (`\`).

For example, if you want to search for a dollar sign in your text, just precede it with a backslash character:

```
$ cat data2
The cost is $4.00
$ sed -n '/\$/p' data2
The cost is $4.00
$
```

Because the backslash is a special character, if you need to use it in a regular expression pattern you'll need to escape it as well, producing a double backslash:

```
$ echo "\ is a special character" | sed -n '\\$/p'
\ is a special character
$
```

Finally, although the forward slash isn't a regular expression special character, if you use it in your regular expression pattern in the sed editor or the gawk program, you'll get an error:

```
$ echo "3 / 2" | sed -n '///p'
sed: -e expression #1, char 2: No previous regular expression
$
```

To use a forward slash you'll need to escape that as well:

```
$ echo "3 / 2" | sed -n '\\//p'
3 / 2
$
```

Now the sed editor can properly interpret the regular expression pattern, and all is well.

Anchor Characters

As shown in the “Plain Text” section, by default, when you specify a regular expression pattern, if the pattern appears anywhere in the data stream, it will match. There are two special characters you can use to anchor a pattern to either the beginning or the end of lines in the data stream.

Starting at the Beginning

The caret character (^) defines a pattern that starts at the beginning of a line of text in the data stream. If the pattern is located any place other than the start of the line of text, the regular expression pattern fails.

To use the caret character, you must place it before the pattern specified in the regular expression:

```
$ echo "The book store" | sed -n '/^book/p'
$
$ echo "Books are great" | sed -n '/^Book/p'
Books are great
$
```

The caret anchor character checks for the pattern at the beginning of each new line of data, as determined by the newline character:

```
$ cat data3
This is a test line.
this is another test line.
A line that tests this feature.
Yet more testing of this
$ sed -n '/^this/p' data3
this is another test line.
$
```

As long as the pattern appears at the start of a new line, the caret anchor will catch it.

If you position the caret character in any place other than at the beginning of the pattern, it will act like a normal character and not as a special character:

```
$ echo "This ^ is a test" | sed -n '/s ^/p'
This ^ is a test
$
```

Because the caret character is listed last in the regular expression pattern, the sed editor uses it as a normal character to match text.

Note

If you need to specify a regular expression pattern using only the caret character, you don't need to escape it with a backslash. However, if you specify the caret character first, followed by additional text in the pattern, you'll need to use the escape character before the caret character.

Looking for the Ending

The opposite of looking for a pattern at the start of a line is looking for it at the end of a line. The dollar sign (\$) special character defines the end anchor. Add this special character after a text pattern to indicate that the line of data must end with the text pattern:

```
$ echo "This is a good book" | sed -n '/book$/p'
This is a good book
$ echo "This book is good" | sed -n '/book$/p'
$
```

The problem with an ending text pattern is that you must

be careful what you're looking for:

```
$ echo "There are a lot of good books" | sed -n '/book$/p'
$
```

Making the word “book” plural at the end of the line means that it no longer matches the regular expression pattern, even though book is in the data stream. The text pattern must be the last thing on the line for the pattern to match.

Combining Anchors

There are a couple of common situations where you can combine both the start and end anchor on the same line. In the first situation, suppose that you want to look for a line of data containing only a specific text pattern:

```
$ cat data4
this is a test of using both anchors
I said this is a test
this is a test
I'm sure this is a test.
$ sed -n '/^this is a test$/p' data4
this is a test
$
```

The sed editor ignores the lines that include other text besides the specified text.

The second situation may seem a little odd at first but is extremely useful. By combining both anchors in a pattern with no text, you can filter blank lines from the data stream. Consider this example:

```
$ cat data5
This is one test line.

This is another test line.
$ sed '/^$/d' data5
```

```
This is one test line.  
This is another test line.  
$
```

The regular expression pattern that is defined looks for lines that have nothing between the start and end of the line. Because blank lines contain no text between the two newline characters, they match the regular expression pattern. The sed editor uses the `d` delete command to delete lines that match the regular expression pattern, thus removing all blank lines from the text. This is an effective way to remove blank lines from documents.

The Dot Character

The dot special character is used to match any single character except a newline character. The dot character must match a character, however; if there's no character in the place of the dot, then the pattern will fail.

Let's take a look at a few examples of using the dot character in a regular expression pattern:

```
$ cat data6  
This is a test of a line.  
The cat is sleeping.  
That is a very nice hat.  
This test is at line four.  
at ten o'clock we'll go home.  
$ sed -n '/.at/p' data6  
The cat is sleeping.  
That is a very nice hat.  
This test is at line four.  
$
```

You should be able to figure out why the first line failed and why the second and third lines passed. The fourth line is a little tricky. Notice that we matched the `at`, but there's no

character in front of it to match the dot character. Ah, but there is! In regular expressions, spaces count as characters, so the space in front of the `at` matches the pattern. The fifth line proves this, by putting the `at` in the front of the line, which fails to match the pattern.

Character Classes

The dot special character is great for matching a character position against any character, but what if you want to limit what characters to match? This is called a *character class* in regular expressions.

You can define a class of characters that would match a position in a text pattern. If one of the characters from the character class is in the data stream, it matches the pattern.

To define a character class, you use square brackets. The brackets should contain any character that you want to include in the class. You then use the entire class within a pattern just like any other wildcard character. This takes a little getting used to at first, but once you catch on it can generate some pretty amazing results.

The following is an example of creating a character class:

```
$ sed -n '/[ch]at/p' data6
The cat is sleeping.
That is a very nice hat.
$
```

Using the same data file as in the dot special character example, we came up with a different result. This time we

managed to filter out the line that just contained the word `at`. The only words that match this pattern are `cat` and `hat`. Also notice that the line that started with `at` didn't match as well. There must be a character in the character class that matches the appropriate position.

Character classes come in handy if you're not sure which case a character is in:

```
$ echo "Yes" | sed -n '/[Yy]es/p'
Yes
$ echo "yes" | sed -n '/[Yy]es/p'
yes
$
```

You can use more than one character class in a single expression:

```
$ echo "Yes" | sed -n '/[Yy][Ee][Ss]/p'
Yes
$ echo "yEs" | sed -n '/[Yy][Ee][Ss]/p'
yEs
$ echo "yeS" | sed -n '/[Yy][Ee][Ss]/p'
yeS
$
```

The regular expression used three character classes to cover both lower and upper cases for all three character positions.

Character classes don't have to contain just letters; you can use numbers in them as well:

```
$ cat data7
This line doesn't contain a number.
This line has 1 number on it.
This line a number 2 on it.
This line has a number 4 on it.
$ sed -n '/[0123]/p' data7
This line has 1 number on it.
This line a number 2 on it.
$
```

The regular expression pattern matches any lines that contain the numbers 0, 1, 2, or 3. Any other numbers are ignored, as are lines without numbers in them.

You can combine character classes to check for properly formatted numbers, such as phone numbers and zip codes. However, when you're trying to match a specific format, you must be careful. Here's an example of a zip code match gone wrong:

```
$ cat data8
60633
46201
223001
4353
22203
$ sed -n `
>/[0123456789][0123456789][0123456789][0123456789][0123456789]/p
>' data8
60633
46201
223001
22203
$
```

This might not have produced the result you were thinking of. It did a fine job of filtering out the number that was too short to be a zip code, as the last character class didn't have a character to match against. However, it still passed the six-digit number, even though we only defined five character classes.

Remember that the regular expression pattern can be found anywhere in the text of the data stream. There can always be additional characters besides the matching pattern characters. If you want to ensure that you only match against five numbers, you need to delineate them somehow, either with spaces, or as in this example, by

showing that they're at the start and end of the line:

```
$ sed -n `
> /^[0123456789][0123456789][0123456789][0123456789][0123456789]$/p
> ` data8
60633
46201
22203
$
```

Now that's much better! Later in this chapter we look at how to simplify this even further.

One extremely popular use for character classes is parsing words that might be misspelled, such as data entered from a user form. You can easily create regular expressions that can accept common misspellings in data:

```
$ cat data9
I need to have some maintenance done on my car.
I'll pay that in a seperate invoice.
After I pay for the maintenance my car will be as good as new.
$ sed -n `
/maint[ea]n[ae]nce/p
/sep[ea]r[ea]te/p
` data9
I need to have some maintenance done on my car.
I'll pay that in a seperate invoice.
After I pay for the maintenance my car will be as good as new.
$
```

The two sed print commands in this example utilize regular expression character classes to help catch the misspelled words, maintenance and separate, in the text. The same regular expression pattern also matches the properly spelled occurrence of “maintenance.”

Negating Character Classes

In regular expression patterns, you can also reverse the

effect of a character class. Instead of looking for a character contained in the class, you can look for any character that's not in the class. To do that, just place a caret character at the beginning of the character class range:

```
$ sed -n '/[^ch]at/p' data6
This test is at line two.
$
```

By negating the character class, the regular expression pattern matches any character that's neither a *c* nor an *h*, along with the text pattern. Because the space character fits this category, it passed the pattern match. However, even with the negation, the character class must still match a character, so the line with the `at` in the start of the line still doesn't match the pattern.

Using Ranges

You may have noticed when I showed the zip code example earlier that it was somewhat awkward having to list all of the possible digits in each character class. Fortunately, you can use a shortcut so you don't have to do that.

You can use a range of characters within a character class by using the dash symbol. Just specify the first character in the range, a dash, and then the last character in the range. The regular expression includes any character that's within the specified character range, according to the character set used by the Linux system (see Chapter 2).

Now you can simplify the zip code example by specifying

a range of digits:

```
$ sed -n '/^[0-9][0-9][0-9][0-9][0-9]$/p' data8
60633
46201
45902
$
```

That saved a lot of typing! Each character class will match any digit from 0 to 9. The pattern will fail if a letter is present anywhere in the data:

```
$ echo "a8392" | sed -n '/^[0-9][0-9][0-9][0-9][0-9]$/p'
$
$ echo "1839a" | sed -n '/^[0-9][0-9][0-9][0-9][0-9]$/p'
$
$ echo "18a92" | sed -n '/^[0-9][0-9][0-9][0-9][0-9]$/p'
$
```

The same technique also works with letters:

```
$ sed -n '/[c-h]at/p' data6
The cat is sleeping.
That is a very nice hat.
$
```

The new pattern `[c-h]at` matches words where the first letter is between the letter *c* and the letter *h*. In this case, the line with only the word `at` failed to match the pattern.

You can also specify multiple, noncontinuous ranges in a single character class:

```
$ sed -n '/[a-ch-m]at/p' data6
The cat is sleeping.
That is a very nice hat.
$
```

The character class allows the ranges *a* through *c*, and *h* through *m* to appear before the `at` text. This range would reject any letters between *d* and *g*:

```
$ echo "I'm getting too fat." | sed -n '/[a-ch-m]at/p'
$
```

This pattern rejected the `fat` text, as it wasn't in the specified range.

Special Character Classes

In addition to defining your own character classes, the BRE contains special character classes you can use to match against specific types of characters. [Table 19.1](#) describes the BRE special characters you can use.

Table 19.1 BRE Special Character Classes

Class	Description
<code>[[[:alpha:]]</code>	Match any alphabetical character, either upper or lower case.
<code>[[[:alnum:]]</code>	Match any alphanumeric character 0–9, A–Z, or a–z
<code>[[[:blank:]]</code>	Match a space or Tab character.
<code>[[[:digit:]]</code>	Match a numerical digit from 0 through 9.
<code>[[[:lower:]]</code>	Match any lowercase alphabetical character a–z.
<code>[[[:print:]]</code>	Match any printable character.
<code>[[[:punct:]]</code>	Match a punctuation character.
<code>[[[:space:]]</code>	Match any whitespace character: space, Tab, NL, FF, VT, CR.
<code>[[[:upper:]]</code>	Match any uppercase alphabetical character A–Z.

You use the special character classes just as you would a normal character class in your regular expression patterns:

```
$ echo "abc" | sed -n '[[[:digit:]]/p'
$
$ echo "abc" | sed -n '[[[:alpha:]]/p'
abc
$ echo "abc123" | sed -n '[[[:digit:]]/p'
abc123
```

```
$ echo "This is, a test" | sed -n '/[[[:punct:]]]/p'
This is, a test
$ echo "This is a test" | sed -n '/[[[:punct:]]]/p'
$
```

Using the special character classes is an easy way to define ranges. Instead of having to use a range [0–9], you can just use `[[[:digit:]]]`.

The Asterisk

Placing an asterisk after a character signifies that the character must appear zero or more times in the text to match the pattern:

```
$ echo "ik" | sed -n '/ie*k/p'
ik
$ echo "iek" | sed -n '/ie*k/p'
iek
$ echo "ieek" | sed -n '/ie*k/p'
ieek
$ echo "ieeek" | sed -n '/ie*k/p'
ieeek
$ echo "ieeeeek" | sed -n '/ie*k/p'
ieeeeek
$
```

This pattern symbol is commonly used for handling words that have a common misspelling or variations in language spellings. For example, if you need to write a script that may be used in either American or British English, you could write:

```
$ echo "I'm getting a color TV" | sed -n '/colou*r/p'
I'm getting a color TV
$ echo "I'm getting a colour TV" | sed -n '/colou*r/p'
I'm getting a colour TV
$
```

The `u*` in the pattern indicates that the letter *u* may or may not appear in the text to match the pattern. Similarly, if you

know of a word that is commonly misspelled, you can accommodate it by using the asterisk:

```
$ echo "I ate a potatoe with my lunch." | sed -n '/potatoe*/p'
I ate a potatoe with my lunch.
$ echo "I ate a potato with my lunch." | sed -n '/potatoe*/p'
I ate a potato with my lunch.
$
```

Placing an asterisk next to the possible extra letter allows you to accept the misspelled word.

Another handy feature is combining the dot special character with the asterisk special character. This combination provides a pattern to match any number of any characters. It's often used between two text strings that may or may not appear next to each other in the data stream:

```
$ echo "this is a regular pattern expression" | sed -n `
> /regular.*expression/p`
this is a regular pattern expression
$
```

Using this pattern, you can easily search for multiple words that may appear anywhere in a line of text in the data stream.

The asterisk can also be applied to a character class. This allows you to specify a group or range of characters that can appear more than once in the text:

```
$ echo "bt" | sed -n '/b[ae]*t/p'
bt
$ echo "bat" | sed -n '/b[ae]*t/p'
bat
$ echo "bet" | sed -n '/b[ae]*t/p'
bet
$ echo "btt" | sed -n '/b[ae]*t/p'
btt
$
$ echo "baat" | sed -n '/b[ae]*t/p'
baat
$ echo "baaaet" | sed -n '/b[ae]*t/p'
```

```
baaeeet
$ echo "baaeaeaat" | sed -n '/b[ae]*t/p'
baaeaeaat
$ echo "baakeeet" | sed -n '/b[ae]*t/p'
$
```

As long as the *a* and *e* characters appear in any combination between the *b* and *t* characters (including not appearing at all), the pattern matches. If any other character outside of the defined character class appears, the pattern match fails.

Extended Regular Expressions

The POSIX ERE patterns include a few additional symbols that are used by some Linux applications and utilities. The *gawk* program recognizes the ERE patterns, but the *sed* editor doesn't.

Caution

It's important to remember that there is a difference between the regular expression engines in the *sed* editor and the *gawk* program. The *gawk* program can use most of the extended regular expression pattern symbols, and it can provide some additional filtering capabilities that the *sed* editor doesn't have. However, because of this, it is often slower in processing data streams.

This section describes the more commonly found ERE pattern symbols that you can use in your *gawk* program scripts.

The Question Mark

The question mark is similar to the asterisk, but with a slight twist. The question mark indicates that the preceding character can appear zero or one time, but that's all. It doesn't match repeating occurrences of the character:

```
$ echo "bt" | gawk '/be?t/{print $0}'
bt
$ echo "bet" | gawk '/be?t/{print $0}'
bet
$ echo "beet" | gawk '/be?t/{print $0}'
$
$ echo "beeeet" | gawk '/be?t/{print $0}'
$
```

If the `e` character doesn't appear in the text, or as long as it appears only once in the text, the pattern matches.

As with the asterisk, you can use the question mark symbol along with a character class:

```
$ echo "bt" | gawk '/b[ae]?t/{print $0}'
bt
$ echo "bat" | gawk '/b[ae]?t/{print $0}'
bat
$ echo "bot" | gawk '/b[ae]?t/{print $0}'
$
$ echo "bet" | gawk '/b[ae]?t/{print $0}'
bet
$ echo "baet" | gawk '/b[ae]?t/{print $0}'
$
$ echo "beat" | gawk '/b[ae]?t/{print $0}'
$
$ echo "beet" | gawk '/b[ae]?t/{print $0}'
$
```

If zero or one character from the character class appears, the pattern match passes. However, if both characters appear, or if one of the characters appears twice, the pattern match fails.

The Plus Sign

The plus sign is another pattern symbol that's similar to the asterisk, but with a different twist than the question mark. The plus sign indicates that the preceding character can appear one or more times, but must be present at least once. The pattern doesn't match if the character is not present:

```
$ echo "beeet" | gawk '/be+t/{print $0}'
beeet
$ echo "beet" | gawk '/be+t/{print $0}'
beet
$ echo "bet" | gawk '/be+t/{print $0}'
bet
$ echo "bt" | gawk '/be+t/{print $0}'
$
```

If the e character is not present, the pattern match fails. The plus sign also works with character classes, the same way as the asterisk and question mark do:

```
$ echo "bt" | gawk '/b[ae]+t/{print $0}'
$
$ echo "bat" | gawk '/b[ae]+t/{print $0}'
bat
$ echo "bet" | gawk '/b[ae]+t/{print $0}'
bet
$ echo "beat" | gawk '/b[ae]+t/{print $0}'
beat
$ echo "beet" | gawk '/b[ae]+t/{print $0}'
beet
$ echo "beeat" | gawk '/b[ae]+t/{print $0}'
beeat
$
```

This time if either character defined in the character class appears, the text matches the specified pattern.

Using Braces

Curly braces are available in ERE to allow you to specify a limit on a repeatable regular expression. This is often referred to as an *interval*. You can express the interval in two formats:

- m - The regular expression appears exactly m times.
- m, n - The regular expression appears at least m times, but no more than n times.

This feature allows you to fine-tune exactly how many times you allow a character (or character class) to appear in a pattern.

Caution

By default, the gawk program doesn't recognize regular expression intervals. You must specify the `--re-interval` command line option for the gawk program to recognize regular expression intervals.

Here's an example of using a simple interval of one value:

```
$ echo "bt" | gawk --re-interval '/be{1}t/{print $0}'
$
$ echo "bet" | gawk --re-interval '/be{1}t/{print $0}'
bet
$ echo "beet" | gawk --re-interval '/be{1}t/{print $0}'
$
```

By specifying an interval of one, you restrict the number of times the character can be present for the string to match the pattern. If the character appears more times, the pattern

match fails.

There are lots of times when specifying the lower and upper limit comes in handy:

```
$ echo "bt" | gawk --re-interval '/be{1,2}t/{print $0}'
$
$ echo "bet" | gawk --re-interval '/be{1,2}t/{print $0}'
bet
$ echo "beet" | gawk --re-interval '/be{1,2}t/{print $0}'
beet
$ echo "beeeet" | gawk --re-interval '/be{1,2}t/{print $0}'
$
```

In this example, the `e` character can appear once or twice for the pattern match to pass; otherwise, the pattern match fails.

The interval pattern match also applies to character classes:

```
$ echo "bt" | gawk --re-interval '/b[ae]{1,2}t/{print $0}'
$
$ echo "bat" | gawk --re-interval '/b[ae]{1,2}t/{print $0}'
bat
$ echo "bet" | gawk --re-interval '/b[ae]{1,2}t/{print $0}'
bet
$ echo "beat" | gawk --re-interval '/b[ae]{1,2}t/{print $0}'
beat
$ echo "beet" | gawk --re-interval '/b[ae]{1,2}t/{print $0}'
beet
$ echo "beeeat" | gawk --re-interval '/b[ae]{1,2}t/{print $0}'
$
$ echo "baeeet" | gawk --re-interval '/b[ae]{1,2}t/{print $0}'
$
$ echo "baeeaat" | gawk --re-interval '/b[ae]{1,2}t/{print $0}'
$
```

This regular expression pattern will match if there are exactly one or two instances of the letter `a` or `e` in the text pattern, but it will fail if there are any more in any combination.

The Pipe Symbol

The pipe symbol allows to you to specify two or more patterns that the regular expression engine uses in a logical OR formula when examining the data stream. If any of the patterns match the data stream text, the text passes. If none of the patterns match, the data stream text fails.

The format for using the pipe symbol is:

```
expr1|expr2|...
```

Here's an example of this:

```
$ echo "The cat is asleep" | gawk '/cat|dog/{print $0}'
The cat is asleep
$ echo "The dog is asleep" | gawk '/cat|dog/{print $0}'
The dog is asleep
$ echo "The sheep is asleep" | gawk '/cat|dog/{print $0}'
$
```

This example looks for the regular expression `cat` or `dog` in the data stream. You can't place any spaces within the regular expressions and the pipe symbol, or they'll be added to the regular expression pattern.

The regular expressions on either side of the pipe symbol can use any regular expression pattern, including character classes, to define the text:

```
$ echo "He has a hat." | gawk '/[ch]at|dog/{print $0}'
He has a hat.
$
```

This example would match `cat`, `hat`, or `dog` in the data stream text.

Grouping Expressions

Regular expression patterns can also be grouped by using parentheses. When you group a regular expression pattern, the group is treated like a standard character. You can apply a special character to the group just as you would to a regular character. For example:

```
$ echo "Sat" | gawk '/Sat(urday)?/{print $0}'
Sat
$ echo "Saturday" | gawk '/Sat(urday)?/{print $0}'
Saturday
$
```

The grouping of the “urday” ending along with the question mark allows the pattern to match either the full day name Saturday or the abbreviated name Sat.

It's common to use grouping along with the pipe symbol to create groups of possible pattern matches:

```
$ echo "cat" | gawk '/(c|b)a(b|t)/{print $0}'
cat
$ echo "cab" | gawk '/(c|b)a(b|t)/{print $0}'
cab
$ echo "bat" | gawk '/(c|b)a(b|t)/{print $0}'
bat
$ echo "bab" | gawk '/(c|b)a(b|t)/{print $0}'
bab
$ echo "tab" | gawk '/(c|b)a(b|t)/{print $0}'
$
$ echo "tac" | gawk '/(c|b)a(b|t)/{print $0}'
$
```

The pattern `(c|b)a(b|t)` matches any combination of the letters in the first group along with any combination of the letters in the second group.

Regular Expressions in Action

Now that you've seen the rules and a few simple

demonstrations of using regular expression patterns, it's time to put that knowledge into action. The following sections demonstrate some common regular expression examples within shell scripts.

Counting Directory Files

To start things out, let's look at a shell script that counts the executable files that are present in the directories defined in your `PATH` environment variable. To do that, you'll need to parse out the `PATH` variable into separate directory names. Chapter 5 showed you how to display the `PATH` environment variable:

```
$ echo $PATH
/usr/local/bin:/bin:/usr/bin:/usr/X11R6/bin:/usr/games:/usr/java/
j2sdk1.4.1_01/bin
$
```

Your `PATH` environment variable will differ, depending on where the applications are located on your Linux system. The key is to recognize that each directory in the `PATH` is separated by a colon. To get a listing of directories that you can use in a script, you'll have to replace each colon with a space. You now recognize that the `sed` editor can do just that using a simple regular expression:

```
$ echo $PATH | sed 's:/ /g'
/usr/local/bin /bin /usr/bin /usr/X11R6/bin /usr/games /usr/java/
j2sdk1.4.1_01/bin
$
```

Once you've got the directories separated out, you can use them in a standard `for` statement (see Chapter 12) to iterate through each directory:

```
mypath='echo $PATH | sed 's:/: /g'`  
for directory in $mypath  
do  
...  
done
```

Once you have each directory, you can use the `ls` command to list each file in each directory, and use another `for` statement to iterate through each file, incrementing a counter for each file.

The final version of the script looks like this:

```
$ cat countfiles  
#!/bin/bash  
# count number of files in your PATH  
mypath='echo $PATH | sed 's:/: /g'`  
count=0  
for directory in $mypath  
do  
    check='ls $directory'  
    for item in $check  
    do  
        count=$(( $count + 1 )  
    done  
    echo "$directory - $count"  
    count=0  
done  
$ ./countfiles  
/usr/local/bin - 79  
/bin - 86  
/usr/bin - 1502  
/usr/X11R6/bin - 175  
/usr/games - 2  
/usr/java/j2sdk1.4.1_01/bin - 27  
$
```

Now we're starting to see some of the power behind regular expressions!

Validating a Phone Number

The previous example showed how to incorporate the

simple regular expression along with `sed` to replace characters in a data stream to process data. Often regular expressions are used to validate data to ensure that data is in the correct format for a script.

A common data validation application checks phone numbers. Often, data entry forms request phone numbers, and often customers fail to enter a properly formatted phone number. In the United States, there are several common ways to display a phone number:

```
(123) 456-7890  
(123) 456-7890  
123-456-7890  
123.456.7890
```

This leaves four possibilities for how customers can enter their phone number in a form. The regular expression must be robust enough to be able to handle either situation.

When building a regular expression, it's best to start on the left-hand side, and build your pattern to match the possible characters you'll run into. In this example, there may or may not be a left parenthesis in the phone number. This can be matched by using the pattern:

```
^\(?
```

The caret is used to indicate the beginning of the data. Because the left parenthesis is a special character, you must escape it to use it as a normal character. The question mark indicates that the left parenthesis may or may not appear in the data to match.

Next comes the three-digit area code. In the United States, area codes start with the number 2 (no area codes

start with the digits 0 or 1), and can go to 9. To match the area code, you'd use the following pattern:

```
[2-9][0-9]{2}
```

This requires that the first character be a digit between 2 and 9, followed by any two digits. After the area code, the ending right parenthesis may or may not be there:

```
\)?
```

After the area code, there can be a space, no space, a dash, or a dot. You can group those using a character group along with the pipe symbol:

```
([ ]-|\.)
```

The very first pipe symbol appears immediately after the left parenthesis to match the no space condition. You must use the escape character for the dot; otherwise, it will be interpreted to match any character.

Next comes the three-digit phone exchange number. Nothing special required here:

```
[0-9]{3}
```

After the phone exchange number, you must match a space, a dash, or a dot (this time you don't have to worry about matching no space because there must be at least a space between the phone exchange number and the rest of the number):

```
([ ]-|\.)
```

Then to finish things off, you must match the four-digit local phone extension at the end of the string:

```
[0-9]{4}$
```

Putting the entire pattern together results in this:

```
^(\?[2-9][0-9]{2}\)?( |1-\.)[0-9]{3}(|1-\.)[0-9]{4}$
```

You can use this regular expression pattern in the `gawk` program to filter out bad phone numbers. All you need to do now is create a simple script using the regular expression in a `gawk` program, and then filter your phone list through the script. Remember that when you use regular expression intervals in the `gawk` program, you must use the `--re-interval` command line option or you won't get the correct results.

Here's the script:

```
$ cat isphone
#!/bin/bash
# script to filter out bad phone numbers
gawk --re-interval '/^\?[2-9][0-9]{2}\)?( |1-\.)[0-9]{3}(|1-\.)[0-9]{4}$/{print $0}'
$
```

While you can't tell from this listing, the `gawk` command is on a single line in the shell script. You can then redirect phone numbers to the script for processing:

```
$ echo "317-555-1234" | ./isphone
317-555-1234
$ echo "000-555-1234" | ./isphone
$
```

Or you can redirect an entire file of phone numbers to filter out the invalid ones:

```
$ cat phonenumber
000-000-0000
123-456-7890
212-555-1234
(317) 555-1234
(202) 555-9876
33523
1234567890
234.123.4567
$ cat phonenumber | ./isphone
212-555-1234
```

(317) 555-1234
(202) 555-9876
234.123.4567
\$

Only the valid phone numbers that match the regular expression pattern appear.

Parsing an E-mail Address

In this day and age, e-mail addresses have become a crucial form of communication. Trying to validate e-mail addresses has become quite a challenge for script builders because there are a myriad of ways to create an e-mail address. The basic form of an e-mail address is:

`username@hostname`

The `username` value can use any alphanumeric character, along with several special characters:

- Dot
- Dash
- Plus sign
- Underscore

These characters can appear in any combination in a valid e-mail userid. The `hostname` portion of the e-mail address consists of one or more domain names and a server name. The server and domain names must also follow strict naming rules, allowing only alphanumeric characters, along with the special characters:

- Dot
- Underscore

The server and domain names are each separated by a dot, with the server name specified first, any subdomain names specified next, and finally, the top-level domain name without a trailing dot.

At one time there were a fairly limited number of top-level domains, and regular expression pattern builders attempted to add them all in patterns for validation. Unfortunately, as the Internet grew so did the possible top-level domains. This technique is no longer a viable solution.

Let's start building the regular expression pattern from the left side. We know that there can be multiple valid characters in the username. This should be fairly easy:

```
^([a-zA-Z0-9_-\.\+])@
```

This grouping specifies the allowable characters in the username, and the plus sign to indicate that there must be at least one character present. The next character is obviously going to be the @ symbol, no surprises there.

The hostname pattern uses the same technique to match the server name and the subdomain names:

```
([a-zA-Z0-9_-\.\+])
```

This pattern matches the text:

```
server
server.subdomain
server.subdomain.subdomain
```

There are special rules for the top-level domain. Top-level domains are only alphabetic characters, and they must be no fewer than two characters (used in country codes) and no more than five characters in length. The following is the regular expression pattern for the top-level

domain:

```
\.([a-zA-Z]{2,5})$
```

Putting the entire pattern together results in the following:

```
^([a-zA-Z0-9_-\.\+])@([a-zA-Z0-9_-\.\+])\.([a-zA-Z]{2,5})$
```

This pattern will filter out poorly formatted e-mail addresses from a data list. Now you can create your script to implement the regular expression:

```
$ echo "rich@here.now" | ./isemail
rich@here.now
$ echo "rich@here.now." | ./isemail
$
$ echo "rich@here.n" | ./isemail
$
$ echo "rich@here-now" | ./isemail
$
$ echo "rich.blum@here.now" | ./isemail
rich.blum@here.now
$ echo "rich_blum@here.now" | ./isemail
rich_blum@here.now
$ echo "rich/blum@here.now" | ./isemail
$
$ echo "rich#blum@here.now" | ./isemail
$
$ echo "rich*blum@here.now" | ./isemail
$
```

Summary

If you manipulate data files in shell scripts, you'll need to become familiar with regular expressions. Regular expressions are implemented in Linux utilities, programming languages, and applications using regular expression engines. A host of different regular expression engines are available in the Linux world. The two most popular are the POSIX Basic Regular Expression (BRE)

engine and the POSIX Extended Regular Expression (ERE) engine. The sed editor conforms mainly to the BRE engine, while the gawk program utilizes most features found in the ERE engine.

A regular expression defines a pattern template that's used to filter text in a data stream. The pattern consists of a combination of standard text characters and special characters. The special characters are used by the regular expression engine to match a series of one or more characters, similarly to how wildcard characters work in other applications.

By combining characters and special characters, you can define a pattern to match most any type of data. You can then use the sed editor or gawk program to filter specific data from a larger data stream, or for validating data received from data entry applications.

The next chapter digs deeper into using the sed editor to perform advanced text manipulation. Lots of advanced features are available in the sed editor that make it useful for handling large data streams and filtering out just what you need.

Chapter 20

Advanced sed

In This Chapter

- Multiline commands
- The hold space
- Negating a command
- Changing the flow
- Pattern replacement
- Using sed in scripts
- Creating sed utilities

Chapter 18 showed you how to use the basics of the sed editor to manipulate text in data streams. The basic sed editor commands are capable of handling most of your everyday text-editing requirements. This chapter takes a look at the more advanced features that the sed editor has to offer. These are features that you might not use as often, but when you need them, it's nice to know that they're there and how to use them.

Multiline Commands

When using the basic sed editor commands, you might have noticed a limitation. All of the sed editor commands perform functions on a single line of data. As the sed editor reads a data stream, it divides the data into lines based on the presence of newline characters. The sed editor processes each line of data one at a time, processing the defined script commands on a line of text, and then moving on to the next line and repeating the process.

There are times when you need to perform actions on data that spans more than one line. This is especially true if you're trying to find or replace a phrase.

For example, if you're looking for the phrase `Linux System Administrators Group` in your data, it's quite possible that the phrase can be split into two lines between any of the words in the phrase. If you processed the text using a normal sed editor command, it would be impossible to detect how the phrase was split.

Fortunately, the designers behind the sed editor thought of that situation and devised a solution. The sed editor includes three special commands that you can use to process multiline text:

- `N`: Add the next line in the data stream to create a multiline group for processing.
- `D`: Delete a single line in a multiline group.
- `P`: Print a single line in a multiline group.

The following sections examine these multiline

commands more closely and demonstrate how you can use them in your scripts.

The next Command

Before you can examine the multiline `next` command, you first need to take a look at how the single-line version of the `next` command works. Once you know what that command does, it's a lot easier to understand how the multiline version of the `next` command operates.

The Single-Line next Command

The lowercase `n` command tells the sed editor to move to the next line of text in the data stream, without going back to the beginning of the commands. Remember that normally the sed editor processes all of the defined commands on a line before moving to the next line of text in the data stream. The single-line `next` command alters this flow.

This may sound somewhat complicated, and sometimes it is. In this example, you have a data file that contains five lines, two of them empty. The goal is to remove the blank line after the header line but leave the blank line before the last line intact. If you write a sed script to just remove blank lines, you will remove both blank lines:

```
$ cat data1
This is the header line.

This is a data line.

This is the last line.
$
```

```
$ sed '/~$/d' datal
This is the header line.
This is a data line.
This is the last line.
$
```

Because the line you want to remove is blank, you don't have any text you can search for to uniquely identify the line. The solution is to use the `n` command. In this next example, the script looks for a unique line, which contains the word `header`. Once the script identifies that line, the `n` command moves the `sed` editor to the next line of text, which is the empty line.

```
$ sed '/header/{n ; d}' datal
This is the header line.
This is a data line.

This is the last line.
$
```

At that point, the `sed` editor continues processing the command list, which uses the `d` command to delete the empty line. When the `sed` editor reaches the end of the command script, it reads the next line of text from the data stream and starts processing commands from the top of the command script. The `sed` editor does not find another line with the word `header`; thus, no further lines are deleted.

Combining Lines of Text

Now that you've seen the single-line `next` command, you can look at the multiline version. The single-line `next` command moves the next line of text from the data stream into the processing space (called the *pattern space*) of the `sed` editor. The multiline version of the `next` command (which

uses a capital `N`) adds the next line of text to the text already in the pattern space.

This has the effect of combining two lines of text from the data stream into the same pattern space. The lines of text are still separated by a newline character, but the sed editor can now treat both lines of text as one line.

Here's a demonstration of how the `N` command operates:

```
$ cat data2
This is the header line.
This is the first data line.
This is the second data line.
This is the last line.
$
$ sed '/first/{ N ; s/\n/ / }' data2
This is the header line.
This is the first data line. This is the second data line.
This is the last line.
$
```

The sed editor script searches for the line of text that contains the word “first” in it. When it finds the line, it uses the `N` command to combine the next line with that line. It then uses the `substitution` command (`s`) to replace the newline character with a space. The result is that the two lines in the text file appear as one line in the sed editor output.

This has a practical application if you're searching for a text phrase that may be split between two lines in the data file. Here's an example:

```
$ cat data3
The first meeting of the Linux System
Administrator's group will be held on Tuesday.
All System Administrators should attend this meeting.
Thank you for your attendance.
$
$ sed 's/System Administrator/Desktop User/' data3
The first meeting of the Linux System
Administrator's group will be held on Tuesday.
```

```
All Desktop Users should attend this meeting.  
Thank you for your attendance.  
$
```

The `substitution` command is looking for the specific two-word phrase `System Administrator` in the text file. In the single line where the phrase appears, everything is fine; the `substitution` command can replace the text. But in the situation where the phrase is split between two lines, the `substitution` command doesn't recognize the matching pattern.

The `N` command helps solve this problem:

```
$ sed 'N ; s/System.Administrator/Desktop User/' data3  
The first meeting of the Linux Desktop User's group will be held on Tuesday.  
All Desktop Users should attend this meeting.  
Thank you for your attendance.  
$
```

By using the `N` command to combine the next line with the line where the first word is found, you can detect when a line split occurs in the phrase.

Notice that the `substitution` command uses a wildcard pattern `(.)` between “System” and “Administration” to match both the space and the newline situation. However, when it matched the newline character, it removed it from the string, causing the two lines to merge into one line. This may not be exactly what you want.

To solve this problem, you can use two `substitution` commands in the `sed` editor script, one to match the multiline occurrence and one to match the single-line occurrence:

```
$ sed '  
> N  
> s/System\nAdministrator/Desktop\nUser/'
```

```
> s/System Administrator/Desktop User/  
> ' data3  
The first meeting of the Linux Desktop  
User's group will be held on Tuesday.  
All Desktop Users should attend this meeting.  
Thank you for your attendance.  
$
```

The first `substitution` command specifically looks for the newline character between the two search words and includes it in the replacement string. This allows you to add the newline character in the same place in the new text.

There's still one subtle problem with this script, however. The script always reads the next line of text into the pattern space before executing the `sed` editor commands. When it reaches the last line of text, there isn't a next line of text to read, so the `N` command causes the `sed` editor to stop. If the matching text is on the last line in the data stream, the commands won't catch the matching data:

```
$ cat data4  
The first meeting of the Linux System  
Administrator's group will be held on Tuesday.  
All System Administrators should attend this meeting.  
$  
$ sed '  
> N  
> s/System\nAdministrator/Desktop\nUser/  
> s/System Administrator/Desktop User/  
> ' data4  
The first meeting of the Linux Desktop  
User's group will be held on Tuesday.  
All System Administrators should attend this meeting.  
$
```

Because the `System Administrator` text appears in the last line in the data stream, the `N` command misses it, as there isn't another line to read into the pattern space to combine. You can easily resolve this problem by moving your single-

line commands before the `N` command and having only the multiline commands appear after the `N` command, like this:

```
$ sed `
> s/System Administrator/Desktop User/
> N
> s/System\nAdministrator/Desktop\nUser/
> ' data4
The first meeting of the Linux Desktop
User's group will be held on Tuesday.
All Desktop Users should attend this meeting.
$
```

Now, the `substitution` command that looks for the phrase in a single line works just fine on the last line in the data stream, and the multiline `substitution` command covers the occurrence in the middle of the data stream.

The Multiline delete Command

Chapter 18 introduced the single-line `delete` command (`d`). The `sed` editor uses it to delete the current line in the pattern space. When working with the `N` command, however, you must be careful when using the single-line `delete` command:

```
$ sed 'N ; /System\nAdministrator/d' data4
All System Administrators should attend this meeting.
$
```

The `delete` command looked for the words `System` and `Administrator` in separate lines and then deleted both of the lines in the pattern space. This may or may not have been what you intended.

The `sed` editor provides the multiline `delete` command (`D`), which deletes only the first line in the pattern space. It removes all characters up to and including the newline

character:

```
$ sed 'N ; /System\nAdministrator/d' data4
Administrator's group will be held on Tuesday.
All System Administrators should attend this meeting.
$
```

The second line of text, added to the pattern space by the `N` command, remains intact. This comes in handy if you need to remove a line of text that appears before a line that you find a data string in.

Here's an example of removing a blank line that appears before the first line in a data stream:

```
$ cat data5

This is the header line.
This is a data line.

This is the last line.
$
$ sed '/^$/N ; /header/D' data5This is the header line.
This is a data line.

This is the last line.
$
```

This `sed` editor script looks for blank lines and then uses the `N` command to add the next line of text into the pattern space. If the new pattern space contents contain the word `header`, the `D` command removes the first line in the pattern space. Without the combination of the `N` and `D` commands, it would be impossible to remove the first blank line without removing all other blank lines.

The Multiline print Command

By now, you're probably catching on to the difference

between the single-line and multiline versions of the commands. The multiline print command (`P`) follows along using the same technique. It prints only the first line in a multiline pattern space. This includes all characters up to the newline character in the pattern space. It is used in much the same way as the single-line `p` command to display text when you use the `-n` option to suppress output from the script.

```
$ sed -n 'N ; /System\nAdministrator/P' data3
The first meeting of the Linux System
$
```

When the multiline match occurs, the `P` command prints only the first line in the pattern space. The power of the multiline `P` command comes into play when you combine it with the `N` and `D` multiline commands.

The `D` command has a unique feature in that it forces the `sed` editor to return to the beginning of the script and repeat the commands on the same pattern space (it doesn't read a new line of text from the data stream). By including the `N` command in the command script, you can effectively single-step through the pattern space, matching multiple lines together.

Next, by using the `P` command, you can print the first line, and then using the `D` command, you can delete the first line and loop back to the beginning of the script. Once you are back at the script's beginning, the `N` command reads in the next line of text and starts the process all over again. This loop continues until you reach the end of the data stream.

The Hold Space

The *pattern space* is an active buffer area that holds the text examined by the sed editor while it processes commands. However, it isn't the only space available in the sed editor for storing text.

The sed editor utilizes another buffer area called the *hold space*. You can use the hold space to temporarily hold lines of text while working on other lines in the pattern space. There are five commands associated with operating with the hold space, as shown in [Table 20.1](#).

Table 20.1 The sed Editor Hold Space Commands

Command	Description
<code>h</code>	Copy pattern space to hold space.
<code>H</code>	Append pattern space to hold space.
<code>g</code>	Copy hold space to pattern space.
<code>G</code>	Append hold space to pattern space.
<code>x</code>	Exchange contents of pattern and hold spaces.

These commands let you copy text from the pattern space to the hold space. This frees up the pattern space to load another string for processing.

Usually, after using the `h` or `H` commands to move a string to the hold space, eventually you want to use the `g`, `G`, or `x` commands to move the stored string back into the pattern space (otherwise, you wouldn't have cared about saving them in the first place).

With two buffer areas, trying to determine what line of text is in which buffer area can sometimes get confusing. Here's a short example that demonstrates how to use the `h` and `g` commands to move data back and forth between the sed editor buffer spaces:

```
$ cat data2
This is the header line.
This is the first data line.
This is the second data line.
This is the last line.
$
$ sed -n '/first/{
> h
> p
> n
> p
> g
> p
> }' data2
This is the first data line.
This is the second data line.
This is the first data line.
$
```

Take a look at the preceding code example step by step:

1. The `sed` script uses a regular expression in the address to filter the line containing the word `first`.
2. When the line containing the word `first` appears, the `h` command places the line in the hold space.
3. The `p` command prints the contents of the pattern space, which is still the first data line.
4. The `n` command retrieves the next line in the data stream (This is the second data line) and places it in the pattern space.
5. The `p` command prints the contents of the pattern space, which is now the second data line.

6. The `g` command places the contents of the hold space (This is the first data line) back into the pattern space, replacing the current text.

7. The `p` command prints the current contents of the pattern space, which is now back to the first data line.

By shuffling the text lines around using the hold space, you are able to force the first data line to appear after the second data line in the output. If you just drop the first `p` command, you can output the two lines in reverse order:

```
$ sed -n '/first/{
> h
> n
> p
> g
> p
> }' data2
This is the second data line.
This is the first data line.
$
```

This is the start of something useful. You can use this technique to create a sed script that reverses an entire file of text lines! To do that, however, you need to see the negating feature of the sed editor, which is what the next section is all about.

Negating a Command

Chapter 18 showed that the sed editor applies commands either to every text line in the data stream or to lines specifically indicated by either a single address or an address range. You can also configure a command to not

apply to a specific address or address range in the data stream.

The exclamation mark command (!) is used to negate a command. This means in situations where the command would normally have been activated, it isn't. Here's an example demonstrating this feature:

```
$ sed -n '/header/!p' data2
This is the first data line.
This is the second data line.
This is the last line.
$
```

The normal `p` command would have printed only the line in the `data2` file that contained the word `header`. By adding the exclamation mark, all of the lines in the file printed except the one that contained the text referenced in the address.

There are several applications where using the exclamation mark comes in handy. Recall that earlier in the chapter, “The next Commands” section showed a situation where a `sed` editor command wouldn't operate on the last line of text in the data stream because there wasn't a line after it. You can use the exclamation point to fix that problem:

```
$ sed 'N; s/System.Administrator/Desktop User/' data4
The first meeting of the Linux Desktop User's group will be held on Tuesday
All System Administrators should attend this meeting.
$
$ sed '$!N; s/System.Administrator/Desktop User/' data4
The first meeting of the Linux Desktop User's group will be held on Tuesday
All Desktop Users should attend this meeting.
$
```

This example shows the exclamation mark used with the `N` command, along with the dollar sign special address. The dollar sign represents the last line of text in the data stream,

so when the sed editor reaches the last line, it doesn't execute the `N` command. For all other lines, however, it does execute the command.

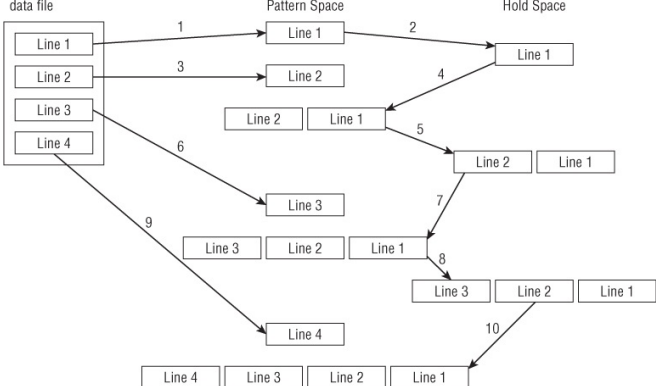
Using this technique, you can reverse the order of text lines in a data stream. To reverse the order of the lines as they appear in the text stream (display the last line first and the first line last), you need to do some fancy footwork using the hold space.

The pattern you'll need to work with goes like this:

1. Place a line in the pattern space.
2. Place the line from the pattern space to the hold space.
3. Put the next line of text in the pattern space.
4. Append the hold space to the pattern space.
5. Place everything in the pattern space into the hold space.
- 6–9. Repeat Steps 3 through 5 until you've put all of the lines in reverse order in the hold space.
10. Retrieve the lines and print them.

[Figure 20.1](#) diagrams what this looks like in more detail.

Figure 20.1 Reversing the order of a text file using the hold space



When using this technique, you do not want to print lines as they are processed. This means using the `-n` command line option for sed. The next thing to determine is how to append the hold space text to the pattern space text. This is done by using the `G` command. The only problem is that you don't want to append the hold space to the first line of text processed. This is easily solved by using the exclamation mark command:

```
1!G
```

The next step is to place the new pattern space (the text line with the appended reverse lines) into the hold space. This is simple enough; just use the `h` command.

When you've got the entire data stream in the pattern space in reverse order, all you need to do is print the results. You know you've got the entire data stream in the

pattern space when you've reached the last line in the data stream. To print the results, just use the following command:

```
$p
```

Those are the pieces you need to create your line-reversing sed editor script. Now try it out in a test run:

```
$ cat data2
This is the header line.
This is the first data line.
This is the second data line.
This is the last line.
$
$ sed -n '{1!G ; h ; $p }' data2
This is the last line.
This is the second data line.
This is the first data line.
This is the header line.
$
```

The sed editor script performed as expected. The output from the script reverses the original lines in the text file. This demonstrates the power of using the hold space in your sed scripts. It provides an easy way to manipulate the order of lines in the script output.

Note

In case you're wondering, there's a Linux command that can perform the function of reversing a text file. The `tac` command displays a text file in reverse order. You probably noticed the clever name of the command because it performs the reverse function of the `cat` command.

Changing the Flow

Normally, the sed editor processes commands starting at the top and proceeding toward the end of the script (the exception is the `D` command, which forces the sed editor to return to the top of the script without reading a new line of text). The sed editor provides a method for altering the flow of the command script, producing a result similar to that of a structured programming environment.

Branching

In the previous section, you saw how the exclamation mark command is used to negate the effect of a command on a line of text. The sed editor provides a way to negate an entire section of commands, based on an address, an address pattern, or an address range. This allows you to perform a group of commands only on a specific subset within the data stream.

The format of the `branch` command is:

```
[address]b [label]
```

The `address` parameter determines which line or lines of data trigger the `branch` command. The `label` parameter defines the location to branch to. If the `label` parameter is not present, the `branch` command proceeds to the end of the script.

```
$ sed '{2,3b ; s/This is/Is this/;s/line./test?/}' data2
Is this the header test?
This is the first data line.
This is the second data line.
```

```
Is this the last test?  
$
```

The `branch` command skips the two `substitution` commands for the second and third lines in the data stream.

Instead of going to the end of the script, you can define a label for the `branch` command to jump to. Labels start with a colon and can be up to seven characters in length:

```
:label2
```

To specify the label, just add it after the `b` command. Using labels allows you to skip commands that match the `branch` address but still process other commands in the script:

```
$ sed '{/first/b jump1 ; s/This is the/No jump on/  
> :jump1  
> s/This is the/Jump here on/}' data2  
No jump on header line  
Jump here on first data line  
No jump on second data line  
No jump on last line  
$
```

The `branch` command specifies that the program should jump to the script line labeled `jump1` if the matching text “first” appears in the line. If the `branch` command pattern doesn't match, the `sed` editor continues processing commands in the script, including the command after the `branch` label. (Thus, all three `substitution` commands are processed on lines that don't match the `branch` pattern.)

If a line matches the `branch` pattern, the `sed` editor branches to the `branch` label line. Thus, only the last `substitution` command is executed.

The example shows branching to a label further down in

the `sed` script. You can also `branch` to a label that appears earlier in the script, thus creating a looping effect:

```
$ echo "This, is, a, test, to, remove, commas. " | sed -n '{
> :start
> s/,//lp
> b start
> }'
This is, a, test, to, remove, commas.
This is a, test, to, remove, commas.
This is a test, to, remove, commas.
This is a test to, remove, commas.
This is a test to remove, commas.
This is a test to remove commas.
```

Each iteration of the script removes the first occurrence of a comma from the text string and prints the string. There's one catch to this script: it never ends. This situation creates an endless loop, searching for commas until you manually stop it by sending a signal with the Ctrl+C key combination.

To prevent this problem, you should specify an address pattern for the `branch` command to look for. If the pattern isn't present, the branching should stop:

```
$ echo "This, is, a, test, to, remove, commas. " | sed -n '{
> :start
> s/,//lp
> /,/b start
> }'
This is, a, test, to, remove, commas.
This is a, test, to, remove, commas.
This is a test, to, remove, commas.
This is a test to, remove, commas.
This is a test to remove, commas.
This is a test to remove commas.
$
```

Now the `branch` command branches only if there's a comma in the line. After the last comma has been removed, the `branch` command won't execute, allowing the script to

properly finish.

Testing

Similar to the `branch` command, the `test` command (`t`) is also used to modify the flow of the `sed` editor script. Instead of jumping to a label based on an address, the `test` command jumps to a label based on the outcome of a `substitution` command.

If the `substitution` command successfully matches and substitutes a pattern, the `test` command branches to the specified label. If the `substitution` command doesn't match the specified pattern, the `test` command doesn't branch.

The `test` command uses the same format as the `branch` command:

```
[address]t [label]
```

Like the `branch` command, if you don't specify a label, `sed` branches to the end of the script if the test succeeds.

The `test` command provides a cheap way to perform a basic `if-then` statement on the text in the data stream. For example, if you don't need to make a substitution if another substitution was made, the `test` command can help:

```
$ sed '{
> s/first/matched/
> t
> s/This is the/No match on/
> }' data2
No match on header line
This is the matched data line
No match on second data line
No match on last line
$
```

The first `substitution` command looks for the pattern text first. If it matches the pattern in the line, it replaces the text, and the `test` command jumps over the subsequent `substitution` command. If the first `substitution` command doesn't match the pattern, the second `substitution` command is processed.

Using the `test` command, you can clean up the loop you tried using the `branch` command:

```
$ echo "This, is, a, test, to, remove, commas. " | sed -n '{
> :start
> s/,//lp
> t start
> }'
This is, a, test, to, remove, commas.
This is a, test, to, remove, commas.
This is a test, to, remove, commas.
This is a test to, remove, commas.
This is a test to remove, commas.
This is a test to remove commas.
$
```

When there are no more substitutions to make, the `test` command doesn't branch and continues with the rest of the script.

Pattern Replacement

You've seen how to use patterns in the `sed` commands to replace text in the data stream. However, when using wildcard characters it's not easy to know exactly what text will match the pattern.

For example, say that you want to place double quotation marks around a word you match in a line. That's simple

enough if you're just looking for one word in the pattern to match:

```
$ echo "The cat sleeps in his hat. " | sed 's/cat/"cat"/'
The "cat" sleeps in his hat.
$
```

But what if you use a wildcard character (.) in the pattern to match more than one word?

```
$ echo "The cat sleeps in his hat. " | sed 's/.at"/.at"/g'
The ".at" sleeps in his ".at".
$
```

The substitution string used the dot wildcard character to match any occurrence of a letter followed by "at". Unfortunately, the replacement string doesn't match the wildcard character value of the matching word.

The Ampersand

The sed editor has a solution for you. The ampersand symbol (&) is used to represent the matching pattern in the substitution command. Whatever text matches the pattern defined, you can use the ampersand symbol to recall it in the replacement pattern. This lets you manipulate whatever word matches the pattern defined:

```
$ echo "The cat sleeps in his hat. " | sed 's/.at/"&"/g'
The "cat" sleeps in his "hat".
$
```

When the pattern matches the word cat, "cat" appears in the substituted word. When it matched the word hat, "hat" appears in the substituted word.

Replacing Individual Words

The ampersand symbol retrieves the entire string that matches the pattern you specify in the `substitution` command. There are times when you'll only want to retrieve a subset of the string. You can do that, too, but it's a little tricky.

The `sed` editor uses parentheses to define a substring component within the substitution pattern. You can then reference each substring component using a special character in the replacement pattern. The replacement character consists of a backslash and a number. The number indicates the position of the substring component. The `sed` editor assigns the first component the character `\1`, the second component the character `\2`, and so on.

Caution

When you use parentheses in the `substitution` command, you must use the escape character to identify them as grouping characters and not normal parentheses. This is the reverse of when you escape other special characters.

Take a look at an example of using this feature in a `sed` editor script:

```
$ echo "The System Administrator manual" | sed `
> s/\(System\) Administrator/\1 User/'
The System User manual
$
```

This `substitution` command uses one set of parentheses

around the word “System,” identifying it as a substring component. It then uses the `\1` in the replacement pattern to recall the first identified component. This wasn't too exciting, but it can really be useful when working with wildcard patterns.

If you need to replace a phrase with just a single word, that's a substring of the phrase, but that substring just happens to be using a wildcard character; using substring components is a lifesaver:

```
$ echo "That furry cat is pretty" | sed 's/furry \(.at\) /\1/'
That cat is pretty
$ echo "That furry hat is pretty" | sed 's/furry \(.at\) /\1/'
That hat is pretty
$
```

In this situation, you can't use the ampersand symbol, as it would replace the entire matching pattern. The substring component provides the answer, allowing you to select just which part of the pattern to use as the replacement pattern.

This feature can be especially helpful when you need to insert text between two or more substring components. Here's a script that uses substring components to insert a comma in long numbers:

```
$ echo "1234567" | sed '{
> :start
> s/\(.*[0-9]\)\([0-9]\{3\}\)/\1,\2/
> t start
> }'
1,234,567
$
```

The script divides the matching pattern into two components:

```
.*[0-9]
[0-9]{3}
```

This pattern looks for two substrings. The first substring is any number of characters, ending in a digit. The second substring is a series of three digits (see Chapter 19 for information about how to use braces in a regular expression). If this pattern is found in the text, the replacement text puts a comma between the two components, each identified by its component position. The script uses the `test` command to iterate through the number until all commas have been placed.

Using sed in Scripts

Now that you've seen the various parts of the sed editor, it's time to put them together and use them in your shell scripts. This section demonstrates some of the features that you should know about when using the sed editor in your bash shell scripts.

Using Wrappers

You may have noticed that trying to implement a sed editor script can be cumbersome, especially if the script is long. Instead of having to retype the entire script each time you want to use it, you can place the sed editor command in a shell script *wrapper*. The wrapper acts as a go-between for the sed editor script and the command line.

Once inside the shell script, you can use normal shell variables and parameters with your sed editor scripts.

Here's an example of using the command line parameter variable as the input to a sed script:

```
$ cat reverse
#!/bin/bash
# shell wrapper for sed editor script to reverse lines

sed -n '{
l!G
h
$p
}' $1
$
```

The shell script called `reverse` uses the sed editor script to reverse text lines in a data stream. It uses the `$1` shell parameter to retrieve the first parameter from the command line, which should be the name of the file to reverse:

```
$ ./reverse data2
This is the last line.
This is the second data line.
This is the first data line.
This is the header line.
$
```

Now you can easily use the sed editor script on any file, without having to constantly retype the entire command line.

Redirecting sed Output

By default, the sed editor outputs the results of the script to `STDOUT`. You can employ all of the standard methods of redirecting the output of the sed editor in your shell scripts.

You can use backticks to redirect the output of your sed editor command to a variable for use later on in the script. The following is an example of using the sed script to add commas to the result of a numeric computation:

```

$ cat fact
#!/bin/bash
#add commas to numbers in factorial answer

factorial=1
counter=1
number=$1

while [ $counter -le $number ]
do
    factorial=$(( $factorial * $counter )
    counter=$(( $counter + 1 )
done

result='echo $factorial | sed '{
:start
s/(.*[0-9])\{([0-9]\{3\}\)\}/\1,\2/
t start
}'

echo "The result is $result"
$
$ ./fact 20
The result is 2,432,902,008,176,640,000
$

```

After you use the normal factorial calculation script, the result of that script is used as the input to the sed editor script, which adds commas. This value is then used in the `echo` statement to produce the result.

Creating sed Utilities

As you've seen in the short examples presented so far in this chapter, there are lots of cool data-formatting things you can do with the sed editor. This section shows a few handy well-known sed editor scripts for performing common data-handling functions.

Double Spacing Lines

To start things off, take a look at a simple sed script to insert a blank line between lines in a text file:

```
$ sed 'G' data2
This is the header line.

This is the first data line.

This is the second data line.

This is the last line.
$
```

That was pretty simple! The key to this trick is the default value of the hold space. Remember that the `G` command simply appends the contents of the hold space to the current pattern space contents. When you start the sed editor, the hold space contains an empty line. By appending that to an existing line, you create a blank line after the existing line.

You may have noticed that this script also adds a blank line to the last line in the data stream, producing a blank line at the end of the file. If you want to get rid of this, you can use the negate symbol and the last line symbol to ensure that the script doesn't add the blank line to the last line of the data stream:

```
$ sed '$!G' data2
This is the header line.

This is the first data line.

This is the second data line.

This is the last line.
$
```

Now that looks a little better. As long as the line isn't the last line, the `G` command appends the contents of the hold space. When the `sed` editor gets to the last line, it skips the `G` command.

Double Spacing Files That May Have Blanks

To take double spacing one step further, what if the text file already has a few blank lines, but you want to double space all of the lines? If you use the previous script, you'll get some areas that have too many blank lines, as each existing blank line gets doubled:

```
$ cat data6
This is line one.
This is line two.

This is line three.
This is line four.
$
$ sed '$!G' data6
This is line one.

This is line two.

This is line three.

This is line four.
$
```

Now you have three blank lines where the original blank line was located. The solution to this problem is to first delete any blank lines from the data stream and then use the `G` command to insert new blank lines after all of the lines. To delete existing blank lines, you just need to use the `d`

command with a pattern that matches a blank line:

```
/$/d
```

This pattern uses the start line tag (the caret) and the end line tag (the dollar sign). Adding this pattern to the script produces the desired results:

```
$ sed '/^$/d;$!G' data6
This is line one.

This is line two.

This is line three.

This is line four.
$
```

Perfect!

Numbering Lines in a File

Chapter 18 showed you how to use the equal sign to display the line numbers of lines in the data stream:

```
$ sed '=' data2
1
This is the header line.
2
This is the first data line.
3
This is the second data line.
4
This is the last line.
$
```

This can be a little awkward to read, as the line number is on a line above the actual line in the data stream. A better solution is to place the line number on the same line as the text.

Now that you've seen how to combine lines using the `N`

command, it shouldn't be too hard to utilize that information in the sed editor script. The trick to this utility, however, is that you can't combine the two commands in the same script.

Once you have the output for the equal sign command, you can pipe the output to another sed editor script that uses the `N` command to combine the two lines. You also need to use the `substitution` command to replace the newline character with either a space or a tab character. Here's what the final solution looks like:

```
$ sed '=' data2 | sed 'N; s/\n/ /'
1 This is the header line.
2 This is the first data line.
3 This is the second data line.
4 This is the last line.
$
```

Now that looks much better. This is a great little utility to have around when working on programs where you need to see the line numbers used in error messages.

Printing Last Lines

So far, you've seen how to use the `p` command to print all of the lines in a data stream or just lines that match a specific pattern. What if you just need to work with the last few lines of a long listing, such as a log file?

The dollar sign represents the last line of a data stream, so it's easy to display just the last line:

```
$ sed -n '$p' data2
This is the last line.
$
```

Now how can you use the dollar sign symbol to display a set number of lines at the end of the data stream? The answer is to create a *rolling window*.

A rolling window is a common way to examine blocks of text lines in the pattern space by combining them using the `N` command. The `N` command appends the next line of text to the text already in the pattern space. Once you have a block of 10 text lines in the pattern space, you can check if you're at the end of the data stream using the dollar sign. If you're not at the end, continue adding more lines to the pattern space, while removing the original lines (remember the `D` command, which deletes the first line in the pattern space).

By looping through the `N` and `D` commands, you add new lines to the block of lines in the pattern space while removing old lines. The `branch` command is the perfect fit for the loop. To end the loop, just identify the last line and use the `q` command to quit.

Here's what the final `sed` editor script looks like:

```
$ sed '{
> :start
> $q
> N
> 11,$D
> b start
> }' /etc/passwd
user:x:1000:1000:user,,,:/home/user:/bin/bash
polkituser:x:113:121:PolicyKit,,,:/var/run/PolicyKit:/bin/false
sshd:x:114:65534:./var/run/sshd:/usr/sbin/nologin
Samantha:x:1001:1002:Samantha,4,,,:/home/Samantha:/bin/bash
Debian-exim:x:115:124:./var/spool/exim4:/bin/false
usbmux:x:116:46:usbmux daemon,,,:/home/usbmux:/bin/false
rtkit:x:117:125:RealtimeKit,,,:/proc:/bin/false
Timothy:x:1002:1005:./home/Timothy:/bin/sh
Christine:x:1003:1006:./home/Christine:/bin/sh
kdm:x:118:65534:./home/kdm:/bin/false
$
```

The script first checks if the line is the last line in the data stream. If it is, the `quit` command stops the loop. The `N` command appends the next line to the current line in the pattern space. The `11,$D` command deletes the first line in the pattern space if the current line is after line 10. This creates the sliding window effect in the pattern space.

Deleting Lines

Another useful utility for the `sed` editor is to remove unwanted blank lines in a data stream. It's easy to remove all the blank lines from a data stream, but it takes a little ingenuity to selectively remove blank lines. This section shows you a couple of quick `sed` editor scripts that you can use to help remove unwanted blank lines from your data.

Deleting Consecutive Blank Lines

It can be a nuisance when extra blank lines crop up in data files. Often you have a data file that contains blank lines, but sometimes a data line is missing and produces too many blank lines (as you saw in the double-spacing example earlier).

The easiest way to remove consecutive blank lines is to check the data stream using a range address. Chapter 18 showed you how to use ranges in addresses, including how to incorporate patterns in the address range. The `sed` editor executes the command for all lines that match within the specified address range.

The key to removing consecutive blank lines is to create an address range that includes a non-blank line and a blank line. If the sed editor comes across this range, it shouldn't delete the line. However, for lines that don't match that range (two or more blank lines in a row), it should delete the lines.

Here's the script to do this:

```
/./,/^$/!d
```

The range is `/./` to `^$/`. The start address in the range matches any line that contains at least one character. The end address in the range matches a blank line. Lines within this range aren't deleted.

Here's the script in action:

```
$ cat data6
This is the first line.

This is the second line.

This is the third line.

This is the fourth line.
$
$ sed '/./,/^$/!d' data6
This is the first line.

This is the second line.

This is the third line.

This is the fourth line.
$
```

No matter how many blank lines appear between lines of data in the file, the output places only one blank line between the lines.

Deleting Leading Blank Lines

It is also a nuisance when data files contain multiple blank lines at the start of the file. Often when trying to import data from a text file into a database, the blank lines create null entries, throwing off any calculations using the data.

Removing blank lines from the top of a data stream is not a difficult task. Here's the script that accomplishes that function:

```
./,,$!d
```

The script uses an address range to determine what lines are deleted. The range starts with a line that contains a character and continues to the end of the data stream. Any line within this range is not deleted from the output. This means that any lines before the first line that contain a character are deleted.

Take a look at this simple script in action:

```
$ cat data7
```

```
This is the first line.
```

```
This is the second line.
```

```
$
```

```
$ sed './,,$!d' data7
```

```
This is the first line.
```

```
This is the second line.
```

```
$
```

The test file contains two blank lines before the data lines. The script successfully removes both of the leading blank lines, while keeping the blank line within the data intact.

Deleting Trailing Blank Lines

Unfortunately, deleting trailing blank lines is not as simple as deleting leading blank lines. Just like printing the end of a data stream, deleting blank lines at the end of a data stream requires a little ingenuity and looping.

Before we start the discussion, let's see what the script looks like:

```
sed '{
: start
/^\\n*$/{$d; N; b start }
}'
```

This may look a little odd to you at first. Notice that there are braces within the normal script braces. This allows you to group commands together within the overall command script. The group of commands applies to the specified address pattern. The address pattern matches any line that contains only a newline character. When one is found, if it's the last line, the `delete` command deletes it. If it's not the last line, the `N` command appends the next line to it, and the `branch` command loops to the beginning to start over.

Here's the script in action:

```
$ cat data8
This is the first line.
This is the second line.

$
$ sed '{
: start
/^\\n*$/{$d; N; b start }
}' data8
This is the first line.
This is the second line.

$
```

The script successfully removed the blank lines from the end of the text file.

Removing HTML Tags

In this day and age, it's not uncommon to download text from a Website to save or use as data in an application. Sometimes, however, when you download text from the Website, you also get the HTML tags used to format the data. This can be a problem when all you want to see is the data.

A standard HTML Web page contains several different types of HTML tags, identifying formatting features required to properly display the page information. Here's a sample of what an HTML file looks like:

```
$ cat data9
<html>
<head>
<title>This is the page title</title>
</head>
<body>
<p>
This is the <b>first</b> line in the Web page. This should provide
some <i>useful</i> information for us to use in our shell script.
</body>
</html>
$
```

HTML tags are identified by the less-than and greater-than symbols. Most HTML tags come in pairs. One tag starts the formatting process (for example, `` for bolding), and another tag stops the formatting process (for example, `` to turn off bolding).

Removing HTML tags creates a problem, however, if

you're not careful. At first glance, you'd think that the way to remove HTML tags would be to just look for a text string that starts with a less-than symbol (<), ends with a greater-than symbol (>), and has data in between the symbols:

```
s/<.*>/g
```

Unfortunately, this command has some unintended consequences:

```
$ sed 's/<.*>/g' data9
```

```
This is the line in the Web page. This should provide
some information for us to use in our shell script.
```

```
$
```

Notice that the title text is missing, along with the text that was bolded and italicized. The sed editor literally interpreted the script to mean any text between the less-than and greater-than sign, including other less-than and greater-than signs! Each time the text was enclosed in HTML tags (such as first), the sed script removed the entire text.

The solution to this problem is to have the sed editor ignore any embedded greater-than signs between the original tags. To do that, you can create a character class that negates the greater-than sign. This changes the script to:

```
s/<[^>]*>/g
```

This script now works properly, displaying the data you need to see from the Web page HTML code:

```
$ sed 's/<[^>]*>/g' data9
```



```
This is the page title
```

```
This is the first line in the Web page. This should provide  
some useful information for us to use in our shell script.
```

```
$
```

That's a little better. To clean things up some, you can add a `delete` command to get rid of those pesky blank lines:

```
$ sed 's/<[^>]*>//g; /^$/d' data9
```

```
This is the page title
```

```
This is the first line in the Web page. This should provide  
some useful information for us to use in our shell script.
```

```
$
```

Now that's much more compact; there's only the data you need to see.

Summary

The `sed` editor provides some advanced features that allow you to work with text patterns across multiple lines. This chapter showed you how to use the `next` command to retrieve the next line in a data stream and place it in the pattern space. Once in the pattern space, you can perform complex `substitution` commands to replace phrases that span more than one line of text.

The `multiline delete` command allows you to remove the first line when the pattern space contains two or more lines. This is a convenient way to iterate through multiple lines in the data stream. Similarly, the `multiline print` command allows you to print just the first line when the pattern space

contains two or more lines of text. The combination of the multiline commands allows you to iterate through the data stream and create a multiline substitution system.

Next, the chapter discussed the hold space. The hold space allows you to set aside a line of text while processing more lines of text. You can recall the contents of the hold space at any time and either replace the text in the pattern space or append the contents of the hold space to the text in the pattern space. Using the hold space allows you to sort through data streams, reversing the order of text lines as they appear in the data.

The chapter also discussed the sed editor flow control commands. The `branch` command provides a way for you to alter the normal flow of sed editor commands in the script, creating loops or skipping commands under certain conditions. The `test` command provides an `if-then` type of statement for your sed editor command scripts. The `test` command branches only if a prior `substitution` command succeeds in replacing text in a line.

The chapter concluded with a discussion of how to use sed scripts in your shell scripts. A common technique for large sed scripts is to place the script in a shell wrapper. You can use command line parameter variables within the sed script to pass shell command line values. This creates an easy way to utilize your sed editor scripts directly from the command line, or even from other shell scripts.

The next chapter digs deeper into the gawk world. The gawk program supports many features of higher-level

programming languages. You can create some pretty involved data manipulation and reporting programs just by using gawk. This chapter will describe the various programming features and demonstrate how to use them to generate your own fancy reports from simple data.

Chapter 21

Advanced gawk

In This Chapter

- Reexamining gawk
- Using variables in gawk
- Using structured commands
- Formatting your printing
- Working with functions

Chapter 18 introduced the gawk program and demonstrated the basics of using it to produce formatted reports from raw data files. This chapter dives more deeply into customizing gawk to produce reports. The gawk program is a full-fledged programming language, providing features that allow you to write advanced programs to manipulate data. If you are jumping into the shell script world from another programming language, you should feel right at home with gawk. In this chapter, you'll see how to use the gawk programming language to write programs to handle just about any data-formatting task you'll run into.

Using Variables

One important feature of any programming language is the ability to store and recall values using variables. The gawk programming language supports two different types of variables:

- Built-in variables
- User-defined variables

Several built-in variables are available for you to use in gawk. The built-in variables contain information used in handling the data fields and records in the data file. You can also create your own variables in your gawk programs. The following sections walk you through how to use variables in your gawk programs.

Built-in Variables

The gawk program uses built-in variables to reference specific features within the program data. This section describes the built-in variables available for you to use in your gawk programs and demonstrates how to use them.

The Field and Record Separator Variables

Chapter 18 demonstrated one type of built-in variable available in gawk, the *data field variables*. The data field variables allow you to reference individual data fields within a data record using a dollar sign and the numerical position

of the data field in the record. Thus, to reference the first data field in the record, you use the `$1` variable. To reference the second data field, you use the `$2` variable, and so on.

Data fields are delineated by a field separator character. By default, the field separator character is a whitespace character, such as a space or a tab. Chapter 18 showed how to change the field separator character either on the command line by using the `-F` command line parameter or within the gawk program by using the special `FS` built-in variable.

The `FS` built-in variable belongs to a group of built-in variables that control how gawk handles fields and records in both input data and output data. [Table 21.1](#) lists the built-in variables contained in this group.

Table 21.1 The gawk Data Field and Record Variables

Variable	Description
<code>FIELDWIDTHS</code>	A space-separated list of numbers defining the exact width (in spaces) of each data field
<code>FS</code>	Input field separator character
<code>RS</code>	Input record separator character
<code>OFS</code>	Output field separator character
<code>ORS</code>	Output record separator character

The `FS` and `OFS` variables define how your gawk program handles data fields in the data stream. You've already seen how to use the `FS` variable to define what character separates data fields in a record. The `OFS` variable performs

the same function but for the output by using the `print` command.

By default, `gawk` sets the `OFS` variable to a space, so when you use the command:

```
print $1,$2,$3
```

you see the output as:

```
field1 field2 field3
```

You can see this in the following example:

```
$ cat data1
data11,data12,data13,data14,data15
data21,data22,data23,data24,data25
data31,data32,data33,data34,data35
$ gawk 'BEGIN{FS=","} {print $1,$2,$3}' data1
data11 data12 data13
data21 data22 data23
data31 data32 data33
$
```

The `print` command automatically places the value of the `OFS` variable between each data field in the output. By setting the `OFS` variable, you can use any string to separate data fields in the output:

```
$ gawk 'BEGIN{FS=","; OFS="-"} {print $1,$2,$3}' data1
data11-data12-data13
data21-data22-data23
data31-data32-data33
$ gawk 'BEGIN{FS=","; OFS="--"} {print $1,$2,$3}' data1
data11--data12--data13
data21--data22--data23
data31--data32--data33
$ gawk 'BEGIN{FS=","; OFS="<-->"} {print $1,$2,$3}' data1
data11<-->data12<-->data13
data21<-->data22<-->data23
data31<-->data32<-->data33
$
```

The `FIELDWIDTHS` variable allows you to read records without using a field separator character. In some applications,

instead of using a field separator character, data is placed in specific columns within the record. In these instances, you must set the `FIELDWIDTHS` variable to match the layout of the data in the records.

Once you set the `FIELDWIDTHS` variable, `gawk` ignores the `FS` and calculates data fields based on the provided field width sizes. Here's an example using field widths instead of field separator characters:

```
$ cat data1b
1005.3247596.37
115-2.349194.00
05810.1298100.1
$ gawk 'BEGIN{FIELDWIDTHS="3 5 2 5"}{print $1,$2,$3,$4}' data1b
100 5.324 75 96.37
115 -2.34 91 94.00
058 10.12 98 100.1
$
```

The `FIELDWIDTHS` variable defines four data fields, and `gawk` parses the data record accordingly. The string of numbers in each record is split based on the defined field width values.

Caution

It's important to remember that once you set the `FIELDWIDTHS` variable, those values must remain constant. This method can't accommodate variable-length data fields.

The `RS` and `ORS` variables define how your `gawk` program handles records in the data stream. By default, `gawk` sets the `RS` and `ORS` variables to the newline character. The default `RS` variable value indicates that each new line of text in the

input data stream is a new record.

Sometimes you run into situations where data fields are spread across multiple lines in the data stream. A classic example of this is data that includes an address and phone number, each on a separate line:

```
Riley Mullen
123 Main Street
Chicago, IL 60601
(312)555-1234
```

If you try to read this data using the default `FS` and `RS` variable values, `gawk` will read each line as a separate record and interpret each space in the record as a field separator. This isn't what you intended.

To solve this problem, you need to set the `FS` variable to the newline character. This indicates that each line in the data stream is a separate field and all of the data on a line belongs to the data field. However, when you do that, you don't know where a new record starts.

To solve this problem, set the `RS` variable to an empty string, and then leave a blank line between data records in the data stream. The `gawk` program will interpret each blank line as a record separator.

The following is an example of using this technique:

```
$ cat data2
Riley Mullen
123 Main Street
Chicago, IL 60601
(312)555-1234

Frank Williams
456 Oak Street
Indianapolis, IN 46201
(317)555-9876
```

```
Haley Snell
4231 Elm Street
Detroit, MI 48201
(313)555-4938
$ gawk 'BEGIN{FS="\n"; RS=""} {print $1,$4}' data2
Riley Mullen (312)555-1234
Frank Williams (317)555-9876
Haley Snell (313)555-4938
$
```

Perfect—the gawk program interpreted each line in the file as a data field and the blank lines as record separators.

Data Variables

Besides the field and record separator variables, gawk provides some other built-in variables to help you know what's going on with your data and extract information from the shell environment. [Table 21.2](#) shows the other built-in variables in gawk.

Table 21.2 More gawk Built-in Variables

Variable	Description
ARGC	The number of command line parameters present.
ARGIND	The index in ARGV of the current file being processed.
ARGV	An array of command line parameters.
CONVFMT	The conversion format for numbers (see the <code>printf</code> statement). The default value is <code>%.6 g</code> .
ENVIRON	An associative array of the current shell environment variables and their values.
ERRNO	The system error if an error occurs when reading or closing input files.
FILENAME	The file name of the data file used for input to the <code>gawk</code> program.
FNR	The current record number in the data file.

IGNORECASE	If set to a non-zero value, ignore the case of characters in strings used in the <code>gawk</code> command.
NF	The total number of data fields in the data file.
NR	The number of input records processed.
OFMT	The output format for displaying numbers. The default is <code>%.6 g</code> .
RLENGTH	The length of the substring matched in the <code>match</code> function.
RSTART	The start index of the substring matched in the <code>match</code> function.

You should recognize a few of these variables from your shell script programming. The `ARGC` and `ARGV` variables allow you to retrieve the number of command line parameters and their values from the shell. This can be a little tricky, however, as `gawk` doesn't count the program script as part of the command line parameters:

```
$ gawk 'BEGIN{print ARGC,ARGV[1]}' data1
2 data1
$
```

The `ARGC` variable indicates that there are two parameters on the command line. This includes the `gawk` command and the `data1` parameter (remember, the program script doesn't count as a parameter). The `ARGV` array starts with an index of 0, which represents the command. The first array value is the first command line parameter after the `gawk` command.

Note

Note that unlike shell variables, when you reference a `gawk` variable in the script, you don't add a dollar sign before the variable name.

The `ENVIRON` variable may seem a little odd to you. It uses an *associative array* to retrieve shell environment variables. An associative array uses text for the array index values instead of numeric values.

The text in the array index is the shell environment variable. The value of the array is the value of the shell environment variable. The following is an example of this:

```
$ gawk '
> BEGIN{
> print ENVIRON["HOME"]
> print ENVIRON["PATH"]
> }'
/home/rich
/usr/local/bin:/bin:/usr/bin:/usr/X11R6/bin
$
```

The `ENVIRON["HOME"]` variable retrieves the `HOME` environment variable value from the shell. Likewise, the `ENVIRON["PATH"]` variable retrieves the `PATH` environment variable value. You can use this technique to retrieve any environment variable value from the shell to use in your gawk programs.

The `FNR`, `NF`, and `NR` variables come in handy when you're trying to keep track of data fields and records in your gawk program. Sometimes you're in a situation where you don't know exactly how many data fields are in a record. The `NF` variable allows you to specify the last data field in the record without having to know its position:

```
$ gawk 'BEGIN{FS=":"; OFS=":"} {print $1,$NF}' /etc/passwd
rich:/bin/bash
testy:/bin/csh
mark:/bin/bash
dan:/bin/bash
mike:/bin/bash
test:/bin/bash
$
```

The `NR` variable contains the numerical value of the last data field in the data file. You can then use it as a data field variable by placing a dollar sign in front of it.

The `FNR` and `NR` variables are similar to each other, but slightly different. The `FNR` variable contains the number of records processed in the current data file. The `NR` variable contains the total number of records processed. Let's look at a couple of examples to see this difference:

```
$ gawk 'BEGIN{FS=","}{print $1,"FNR="FNR}' data1 data1
data11 FNR=1
data21 FNR=2
data31 FNR=3
data11 FNR=1
data21 FNR=2
data31 FNR=3
$
```

In this example, the gawk program command line defines two input files. (It specifies the same input file twice.) The script prints the first data field value and the current value of the `FNR` variable. Notice that the `FNR` value was reset back to 1 when the gawk program processed the second data file.

Now, let's add the `NR` variable and see what that produces:

```
$ gawk `
> BEGIN {FS=","}
> {print $1,"FNR="FNR,"NR="NR}
> END{print "There were",NR,"records processed"}' data1 data1
data11 FNR=1 NR=1
data21 FNR=2 NR=2
data31 FNR=3 NR=3
data11 FNR=1 NR=4
data21 FNR=2 NR=5
data31 FNR=3 NR=6
There were 6 records processed
$
```

The `FNR` variable value was reset when gawk processed

the second data file, but the `NR` variable maintained its count into the second data file. The bottom line is that if you're using only one data file for input, the `FNR` and `NR` values will be the same. If you're using multiple data files for input, the `FNR` value will be reset for each data file, and the `NR` value will keep count throughout all the data files.

Note

You'll notice when using `gawk` that often the `gawk` script can become larger than the rest of your shell script. In the examples in this chapter, for simplicity we just run the `gawk` scripts directly from the command line, using the multiline feature of the shell. When you use `gawk` in a shell script, you should place different `gawk` commands on separate lines. This will make it much easier for you to read and follow, rather than trying to cram it all onto one line in the shell script. Also, if you find yourself using the same `gawk` scripts in different shell scripts, remember you can save the `gawk` script in a separate file and reference it using the `-f` parameter (see Chapter 18).

User-Defined Variables

Just like any other self-respecting programming language, `gawk` allows you to define your own variables for use within the program code. A `gawk` user-defined variable name can be any number of letters, digits, and underscores, but it can't begin with a digit. It is also important to remember that `gawk` variable names are case sensitive.

Assigning Variables in Scripts

Assigning values to variables in gawk programs is similar to doing so in a shell script, using an *assignment statement*:

```
$ gawk `
> BEGIN{
> testing="This is a test"
> print testing
> }'
This is a test
$
```

The output of the `print` statement is the current value of the `testing` variable. Like shell script variables, gawk variables can hold either numeric or text values:

```
$ gawk `
> BEGIN{
> testing="This is a test"
> print testing
> testing=45
> print testing
> }'
This is a test
45
$
```

In this example, the value of the `testing` variable is changed from a text value to a numeric value.

Assignment statements can also include mathematical algorithms to handle numeric values:

```
$ gawk 'BEGIN{x=4; x= x * 2 + 3; print x}'
11
$
```

As you can see from this example, the gawk programming language includes the standard mathematical operators for processing numerical values.

These can include the remainder symbol (%) and the exponentiation symbol (using either ^ or **).

Assigning Variables on the Command Line

You can also use the gawk command line to assign values to variables for the gawk program. This allows you to set values outside of the normal code, changing values on the fly. Here's an example of using a command line variable to display a specific data field in the file:

```
$ cat script1
BEGIN{FS=","}
{print $n}
$ gawk -f script1 n=2 data1
data12
data22
data32
$ gawk -f script1 n=3 data1
data13
data23
data33
$
```

This feature allows you to change the behavior of the script without necessitating that you change the actual script code. The first example displays the second data field in the file, while the second example displays the third data field, just by setting the value of the `n` variable in the command line.

There's one problem with using command line parameters to define variable values. When you set the variable, the value isn't available in the `BEGIN` section of the code:


```
$ cat script2
BEGIN{print "The starting value is",n; FS=","}
{print $n}
$ gawk -f script2 n=3 data1
The starting value is
data13
data23
data33
$
```

You can solve this using the `-v` command line parameter. This allows you to specify variables that are set before the `BEGIN` section of code. The `-v` command line parameter must be placed before the script code in the command line:

```
$ gawk -v n=3 -f script2 data1
The starting value is 3
data13
data23
data33
$
```

Now the `n` variable contains the value set in the command line during the `BEGIN` section of code.

Working with Arrays

Many programming languages provide arrays for storing multiple values in a single variable. The gawk programming language provides the array feature using *associative arrays*.

Associative arrays are different from numerical arrays in that the index value can be any text string. You don't have to use sequential numbers to identify data elements contained in the array. Instead, an associative array consists of a hodge-podge of strings referencing values. Each index

string must be unique and uniquely identifies the data element that's assigned to it. If you're familiar with other programming languages, this is the same concept as hash maps or dictionaries.

The following sections walk you through using associative array variables in your gawk programs.

Defining Array Variables

You can define an array variable using a standard assignment statement. The format of the array variable assignment is:

```
var[index] = element
```

where *var* is the variable name, *index* is the associative array index value, and *element* is the data element value. Here are some examples of array variables in gawk:

```
capital["Illinois"] = "Springfield"  
capital["Indiana"] = "Indianapolis"  
capital["Ohio"] = "Columbus"
```

When you reference an array variable, you must include the index value to retrieve the appropriate data element value:

```
$ gawk 'BEGIN{  
> capital["Illinois"] = "Springfield"  
> print capital["Illinois"]  
> }'  
Springfield  
$
```

When you reference the array variable, the data element value appears. This also works with numeric data element values:

```
$ gawk 'BEGIN{
> var[1] = 34
> var[2] = 3
> total = var[1] + var[2]
> print total
> }'
37
$
```

As you can see from this example, you can use array variables just as you would any other variable in the gawk program.

Iterating Through Array Variables

The problem with associative array variables is that you might not have any way of knowing what the index values are. Unlike numeric arrays, which use sequential numbers for index values, an associative array index can be anything.

If you need to iterate through an associate array in gawk, you can use a special format of the `for` statement:

```
for (var in array)
{
    statements
}
```

The `for` statement loops through the statements, each time assigning the variable `var` the next index value from the `array` associative array. It's important to remember that the variable is the value of the index and not the data element value. You can easily extract the data element value by using the variable as the array index:

```
$ gawk 'BEGIN{
> var["a"] = 1
> var["g"] = 2
```

```

$ var["m"] = 3
> var["u"] = 4
> for (test in var)
> {
>     print "Index:",test," - Value:",var[test]
> }
> }'
Index: u - Value: 4
Index: m - Value: 3
Index: a - Value: 1
Index: g - Value: 2
$

```

Notice that the index values aren't returned in any particular order, but they each reference the appropriate data element value. This is somewhat important to know, as you can't count on the returned values being in the same order, just that the index and data values match.

Deleting Array Variables

Removing an array index from an associative array requires a special command:

```
delete array[index]
```

The `delete` command removes the associative index value and the associated data element value from the array:

```

$ gawk 'BEGIN{
> var["a"] = 1
> var["g"] = 2
> for (test in var)
> {
>     print "Index:",test," - Value:",var[test]
> }
> delete var["g"]
> print "---"
> for (test in var)
>     print "Index:",test," - Value:",var[test]
> }'
Index: a - Value: 1
Index: g - Value: 2

```

```
--  
Index: a - Value: 1  
$
```

Once you delete an index value from the associative array, you can't retrieve it.

Using Patterns

The gawk program supports several types of matching patterns to filter data records, in much the same way as the sed editor. Chapter 18 already showed two special patterns in action. The `BEGIN` and `END` keywords are special patterns that execute statements before or after the data stream data has been read. Similarly, you can create other patterns to execute statements when matching data appears in the data stream.

This section demonstrates how to use matching patterns in your gawk scripts to limit what records a program script applies to.

Regular Expressions

Chapter 19 showed how to use regular expressions as matching patterns. You can use either a Basic Regular Expression (BRE) or an Extended Regular Expression (ERE) to filter which lines in the data stream the program script applies to.

When using a regular expression, the regular expression must appear before the left brace of the program script that

it controls:

```
$ gawk 'BEGIN{FS=","} /11/{print $1}' data1
data11
$
```

The regular expression `/11/` matches records that contain the string `11` anywhere in the data fields. The `gawk` program matches the defined regular expression against all the data fields in the record, including the field separator character:

```
$ gawk 'BEGIN{FS=","} /,d/{print $1}' data1
data11
data21
data31
$
```

This example matches the comma used as the field separator in the regular expression. This is not always a good thing. It can lead to problems trying to match data specific to one data field that may also appear in another data field. If you need to match a regular expression to a specific data instance, you should use the matching operator.

The Matching Operator

The *matching operator* allows you to restrict a regular expression to a specific data field in the records. The matching operator is the tilde symbol (`~`). You specify the matching operator, along with the data field variable, and the regular expression to match:

```
$1 ~ /^data/
```

The `$1` variable represents the first data field in the record. This expression filters records where the first data

field starts with the text data. The following is an example of using the matching operator in a gawk program script:

```
$ gawk 'BEGIN{FS=","} $2 ~ /^data2/{print $0}' data1
data21,data22,data23,data24,data25
$
```

The matching operator compares the second data field with the regular expression `/^data2/`, which indicates the string starts with the text `data2`.

This is a powerful tool that is commonly used in gawk program scripts to search for specific data elements in a data file:

```
$ gawk -F: '$1 ~ /rich/{print $1,$NF}' /etc/passwd
rich /bin/bash
$
```

This example searches the first data field for the text `rich`. When it finds the pattern in a record, it prints the first and last data field values of the record.

You can also negate the regular expression match by using the `!` symbol:

```
$1 !~ /expression/
```

If the regular expression isn't found in the record, the program script is applied to the record data:

```
$ gawk -F: '$1 !~ /rich/{print $1,$NF}' /etc/passwd
root /bin/bash
daemon /bin/sh
bin /bin/sh
sys /bin/sh
--- output truncated ---
$
```

In this example, the gawk program script prints the userid and shell for all of the entries in the `/etc/passwd` file that don't match the userid `rich`!

Mathematical Expressions

In addition to regular expressions, you can also use mathematical expressions in the matching pattern. This feature comes in handy when matching numerical values in data fields. For example, if you want to display all of the system users who belong to the root users group (group number 0), you could use this script:

```
$ gawk -F: '$4 == 0{print $1}' /etc/passwd
root
sync
shutdown
halt
operator
$
```

The script checks for records where the fourth data field contains the value 0. On this Linux system there are five user accounts that belong to the root user group.

You can use any of the normal mathematical comparison expressions:

- $x == y$: Value x is equal to y .
- $x <= y$: Value x is less than or equal to y .
- $x < y$: Value x is less than y .
- $x >= y$: Value x is greater than or equal to y .
- $x > y$: Value x is greater than y .

You can also use expressions with text data, but you must be careful. Unlike regular expressions, expressions are an exact match. The data must match exactly with the pattern:

```
$ gawk -F, '$1 == "data"{print $1}' data1
```



```
$ gawk -F, '$1 == "data11" {print $1}' data1
data11
$
```

The first test doesn't match any records as the first data field value isn't data in any of the records. The second test matches one record with the value `data11`.

Structured Commands

The gawk programming language supports the usual cast of structured programming commands. This section describes each of these commands and demonstrates how to use them within a gawk programming environment.

The if Statement

The gawk programming language supports the standard `if-then-else` format of the `if` statement. You must define a condition for the `if` statement to evaluate, enclosed in parentheses. If the condition evaluates to a `TRUE` condition, the statement immediately following the `if` statement is executed. If the condition evaluates to a `FALSE` condition, the statement is skipped. This can use the format:

```
if (condition)
    statement1
```

or you can place it on one line, like this:

```
if (condition) statement1
```

Here's a simple example demonstrating this format:

```
$ cat data4
```

```
10
5
13
50
34
$ gawk '{if ($1 > 20) print $1}' data4
50
34
$
```

Not too complicated. If you need to execute multiple statements in the `if` statement, you must enclose them with braces:

```
$ gawk '{
> if ($1 > 20)
> {
>   x = $1 * 2
>   print x
> }
> }' data4
100
68
$
```

Be careful that you don't confuse the `if` statement braces with the braces used to start and stop the program script. The `gawk` program can detect missing braces and will produce an error message if you mess up:

```
$ gawk '{
> if ($1 > 20)
> {
>   x = $1 * 2
>   print x
> }' data4
gawk: cmd. line:7: (END OF FILE)
gawk: cmd. line:7: parse error
$
```

The `gawk if` statement also supports the `else` clause, allowing you to execute one or more statements if the `if` statement condition fails. Here's an example of using the `else` clause:

```
$ gawk '{
> if ($1 > 20)
> {
>     x = $1 * 2
>     print x
> } else
> {
>     x = $1 / 2
>     print x
> } }' data4
5
2.5
6.5
100
68
$
```

You can use the `else` clause on a single line, but you must use a semicolon after the `if` statement section:

```
if (condition) statement1; else statement2
```

Here's the same example using the single line format:

```
$ gawk '{if ($1 > 20) print $1 * 2; else print $1 / 2}' data4
5
2.5
6.5
100
68
$
```

This format is more compact but can be harder to follow.

The while Statement

The `while` statement provides a basic looping feature for gawk programs. The following is the format of the `while` statement:

```
while (condition)
{
    statements
}
```

The `while` loop allows you to iterate over a set of data, checking a condition that stops the iteration. This is useful if you have multiple data values in each record that you must use in calculations:

```
$ cat data5
130 120 135
160 113 140
145 170 215
$ gawk '{
> total = 0
> i = 1
> while (i < 4)
> {
>     total += $i
>     i++
> }
> avg = total / 3
> print "Average:", avg
> }' data5
Average: 128.333
Average: 137.667
Average: 176.667
$
```

The `while` statement iterates through the data fields in the record, adding each value to the total variable, and then incrementing the counter variable, `i`. When the counter value is equal to 4, the `while` condition becomes `FALSE`, and the loop terminates, dropping through to the next statement in the script. That statement calculates the average; then the average is printed. This process is repeated for each record in the data file.

The `gawk` programming language supports using the `break` and `continue` statements in `while` loops, allowing you to jump out of the middle of the loop:

```
$ gawk '{
> total = 0
> i = 1
```

```

$ while (i < 4)
> {
>     total += $i
>     if (i == 2)
>         break
>     i++
> }
> avg = total / 2
> print "The average of the first two data elements is:",avg
> }' data5
The average of the first two data elements is: 125
The average of the first two data elements is: 136.5
The average of the first two data elements is: 157.5
$

```

The `break` statement is used to break out of the `while` loop if the value of the `i` variable is 2.

The do-while Statement

The `do-while` statement is similar to the `while` statement but performs the statements before checking the condition statement. The following is the format for the `do-while` statement:

```

do
{
    statements
} while (condition)

```

This format guarantees that the statements are executed at least one time before the condition is evaluated. This comes in handy when you need to perform statements before evaluating the condition:

```

$ gawk '{
> total = 0
> i = 1
> do
> {
>     total += $i
>     i++
> } while (total < 150)

```

```
> print total }' data5
250
160
315
$
```

The script reads the data fields from each record and totals them until the cumulative value reaches 150. If the first data field is over 150 (as seen in the second record), the script is guaranteed to read at least the first data field before evaluating the condition.

The for Statement

The `for` statement is a common method using in many programming languages for looping. The gawk programming language supports the C-style of `for` loops:

```
for( variable assignment; condition; iteration process)
```

This helps simplify the loop by combining several functions in one statement:

```
$ gawk '{
> total = 0
> for (i = 1; i < 4; i++)
> {
>     total += $i
> }
> avg = total / 3
> print "Average:",avg
> }' data5
Average: 128.333
Average: 137.667
Average: 176.667
$
```

By defining the iteration counter in the `for` loop, you don't have to worry about incrementing it yourself as you did when using the `while` statement.

Formatted Printing

You may have noticed that the `print` statement doesn't exactly give you much control over how `gawk` displays your data. About all you can do is control the output field separator character (`OFS`). If you're creating detailed reports, often you'll need to place data in a specific format and location.

The solution is to use the formatted printing command, called `printf`. If you're familiar with C programming, the `printf` command in `gawk` performs the same way, allowing you to specify detailed instructions on how to display data.

The following is the format of the `printf` command:

```
printf "format string", var1, var2 . . .
```

The *format string* is the key to the formatted output. It specifies exactly how the formatted output should appear, using both text elements and *format specifiers*. A format specifier is a special code that indicates what type of variable is displayed and how to display it. The `gawk` program uses each format specifier as a placeholder for each variable listed in the command. The first format specifier matches the first variable listed, the second matches the second variable, and so on.

The format specifiers use the following format:

```
%[modifier]control-letter
```

where *control-letter* is a one-character code that indicates what type of data value will be displayed, and *modifier*

defines an optional formatting feature.

[Table 21.3](#) lists the control letters that can be used in the format specifier.

Table 21.3 Format Specifier Control Letters

Control Letter	Description
c	Displays a number as an ASCII character.
d	Displays an integer value.
i	Displays an integer value (same as <code>d</code>).
e	Displays a number in scientific notation.
f	Displays a floating-point value.
g	Displays either scientific notation or floating point, whichever is shorter.
o	Displays an octal value.
s	Displays a text string.
x	Displays a hexadecimal value.
X	Displays a hexadecimal value, but using capital letters for <i>A</i> through <i>F</i> .

Thus, if you need to display a string variable, you'd use the format specifier `%s`. If you need to display an integer variable, you'd use either `%d` or `%i` (`%d` is the C-style for decimals). If you want to display a large value using scientific notation, you'd use the `%e` format specifier:

```
$ gawk 'BEGIN{
> x = 10 * 100
> printf "The answer is: %e\n", x
> }'
The answer is: 1.000000e+03
$
```


In addition to the control letters, there are three modifiers that you can use for even more control over your output:

- **width:** A numeric value that specifies the minimum width of the output field. If the output is shorter, `printf` pads the space with spaces, using right justification for the text. If the output is longer than the specified width, it overrides the `width` value.
- **prec:** A numeric value that specifies the number of digits to the right of the decimal place in floating-point numbers, or the maximum number of characters displayed in a text string.
- **- (minus sign):** The minus sign indicates that left justification should be used instead of right justification when placing data in the formatted space.

When using the `printf` statement, you have complete control over how your output appears. For example, in the “Built-in Variables” section, we used the `print` command to display data fields from our records:

```
$ gawk 'BEGIN{FS="\n"; RS=""} {print $1,$4}' data2
Riley Mullen (312)555-1234
Frank Williams (317)555-9876
Haley Snell (313)555-4938
$
```

You can use the `printf` command to help format the output so it looks better. First, let's just convert the `print` command to a `printf` command and see what that does:

```
$ gawk 'BEGIN{FS="\n"; RS=""} {printf "%s %s\n", $1, $4}' data 2
Riley Mullen (312)555-1234
Frank Williams (317)555-9876
Haley Snell (313)555-4938
$
```

That produces the same output as the `print` command. The `printf` command uses the `%s` format specifier as a placeholder for the two string values.

Notice that you have to manually add the newline character at the end of the `printf` command to force a new line. Without it, the `printf` command will continue to use the same line on subsequent prints.

This is useful if you need to print multiple things on the same line, but using separate `printf` commands:

```
$ gawk 'BEGIN{FS=","} {printf "%s ", $1} END{printf "\n"}' data1
data11 data21 data31
$
```

Each of the `printf` outputs appears on the same line. To be able to terminate the line, the `END` section prints a single newline character.

Next, let's use a modifier to format the first string value:

```
$ gawk 'BEGIN{FS="\n"; RS=""} {printf "%16s  %s\n", $1, $4}' data2
Riley Mullen      (312)555-1234
Frank Williams    (317)555-9876
Haley Snell       (313)555-4938
$
```

By adding the `16` modifier value, we force the output for the first string to use 16 spaces. By default, the `printf` command uses right justification to place the data in the format space. To make it left justified, just add a minus sign to the modifier:

```
$ gawk 'BEGIN{FS="\n"; RS=""} {printf "%-16s  %s\n", $1, $4}' data2
Riley Mullen      (312)555-1234
Frank Williams    (317)555-9876
Haley Snell       (313)555-4938
$
```

Now that looks pretty professional!

The `printf` command also comes in handy when dealing with floating-point values. By specifying a format for the variable, you can make the output look more uniform:

```
$ gawk '{
> total = 0
> for (i = 1; i < 4; i++)
> {
>     total += $i
> }
> avg = total / 3
> printf "Average: %5.1f\n",avg
> }' data5
Average: 128.3
Average: 137.7
Average: 176.7
$
```

By using the `%5.1f` format specifier, you can force the `printf` command to round the floating-point values to a single decimal place.

Built-in Functions

The gawk programming language provides quite a few built-in functions that perform common mathematical, string, and even time functions. You can utilize these functions in your gawk programs to help cut down on the coding requirements in your scripts. This section walks you through the different built-in functions available in gawk.

Mathematical Functions

If you've done programming in any type of language, you're probably familiar with using built-in functions in your code to

perform common mathematical functions. The gawk programming language doesn't disappoint those looking for advanced mathematical features.

[Table 21.4](#) shows the mathematical built-in functions available in gawk.

[Table 21.4](#) The gawk Mathematical Functions

Function	Description
<code>atan2(x, y)</code>	The arctangent of x / y , with x and y specified in radians.
<code>cos(x)</code>	The cosine of x , with x specified in radians.
<code>exp(x)</code>	The exponential of x .
<code>int(x)</code>	The integer part of x , truncated toward 0.
<code>log(x)</code>	The natural logarithm of x .
<code>rand()</code>	A random floating point value larger than 0 and less than 1.
<code>sin(x)</code>	The sine of x , with x specified in radians.
<code>sqrt(x)</code>	The square root of x .
<code>srand(x)</code>	Specify a seed value for calculating random numbers.

While it does not have an extensive list of mathematical functions, gawk does provide some of the basic elements you need for standard mathematical processing. The `int()` function produces the integer portion of a value, but it doesn't round the value. It behaves much like a floor function found in other programming languages. It produces the nearest integer to a value between the value and 0.

This means that the `int()` function of the value 5.6 will return 5, while the `int()` function of the value -5.6 will return -5.

The `rand()` function is great for creating random numbers, but you'll need to use a trick to get meaningful values. The `rand()` function returns a random number, but only between the values 0 and 1 (not including 0 or 1). To get a larger number, you'll need to scale the returned value.

A common method for producing larger integer random numbers is to create an algorithm that uses the `rand()` function, along with the `int()` function:

```
x = int(10 * rand())
```

This returns a random integer value between (and including) 0 and 9. Just substitute the 10 in the equation with the upper limit value for your application, and you're ready to go.

Be careful when using some of the mathematical functions, as the gawk programming language does have a limited range of numeric values it can work with. If you go over that range, you'll get an error message:

```
$ gawk 'BEGIN{x=exp(100); print x}'
26881171418161356094253400435962903554686976
$ gawk 'BEGIN{x=exp(1000); print x}'
gawk: warning: exp argument 1000 is out of range
inf
$
```

The first example calculates the exponential of 100, which is a very large number but within the range of the system. The second example attempts to calculate the exponential of 1000, which goes over the numerical range limit of the system and produces an error message.

Besides the standard mathematical functions, gawk also provides a few functions for bitwise manipulating of data:

- `and(v1, v2)`: Performs a bitwise **AND** of values `v1` and `v2`.
- `compl(val)`: Performs the bitwise complement of `val`.
- `lshift(val, count)`: Shifts the value `val` count number of bits left.
- `or(v1, v2)`: Performs a bitwise **OR** of values `v1` and `v2`.
- `rshift(val, count)`: Shifts the value `val` count number of bits right.
- `xor(v1, v2)`: Performs a bitwise **XOR** of values `v1` and `v2`.

The bit manipulation functions are useful when working with binary values in your data.

String Functions

The gawk programming language also provides several functions you can use to manipulate string values, shown in [Table 21.5](#).

Table 21.5 The gawk String Functions

Function	Description
<code>asort(s [,d])</code>	Sort an array <code>s</code> based on the data element values. The index values are replaced with sequential numbers indicating the new sort order. Alternatively, the new sorted array is stored in array <code>d</code> if specified.
<code>asorti(s [,d])</code>	Sort an array <code>s</code> based on the index values. The resulting array contains the index values as the data element values, with sequential number indexes indicating the sort order. Alternatively, the new sorted array is stored in array <code>d</code> if specified.
<code>gensub(r, s, h)</code>	Search either the variable <code>s</code> , or the target string <code>t</code> if supplied, for matches of the regular expression <code>r</code> . If <code>h</code> is a string beginning with either <code>o</code> or <code>r</code> , replace the matching

<code>[, t])</code>	string beginning with either <code>g</code> or <code>G</code> , replaces the matching text with <code>s</code> . If <code>h</code> is a number, it represents which occurrence of <code>r</code> to replace.
<code>gsub(r, s [,t])</code>	Search either the variable <code>\$0</code> , or the target string <code>t</code> if supplied, for matches of the regular expression <code>r</code> . If found, substitute the string <code>s</code> globally.
<code>index(s, t)</code>	Returns the index of the string <code>t</code> in string <code>s</code> , or 0 if not found.
<code>length([s])</code>	Returns the length of string <code>s</code> , or if not specified, the length of <code>\$0</code> .
<code>match(s, r [,a])</code>	Returns the index of the string <code>s</code> where the regular expression <code>r</code> occurs. If array <code>a</code> is specified, it contains the portion of <code>s</code> that matches the regular expression.
<code>split(s, a [,r])</code>	Splits <code>s</code> into array <code>a</code> using either the FS character, or the regular expression <code>r</code> if supplied. Returns the number of fields.
<code>sprintf(format, variables)</code>	Returns a string similar to the output of <code>printf</code> using the <code>format</code> and <code>variables</code> supplied.
<code>sub(r, s [,t])</code>	Search either the variable <code>\$0</code> , or the target string <code>t</code> , for matches of the regular expression <code>r</code> . If found, substitutes the string <code>s</code> for the first occurrence.
<code>substr(s, i [,n])</code>	Returns the <code>n</code> th character substring of <code>s</code> , starting at index <code>i</code> . If <code>n</code> is not supplied, the rest of <code>s</code> is used.
<code>tolower(s)</code>	Converts all characters in <code>s</code> to lowercase.
<code>toupper(s)</code>	Converts all characters in <code>s</code> to uppercase.

Some of the string functions are relatively self-explanatory:

```
$ gawk 'BEGIN{x = "testing"; print toupper(x); print length(x) }'
TESTING
7
$
```

However, some of the string functions can get pretty

complicated. The `asort` and `asorti` functions are new gawk functions that allow you to sort an array variable based on either the data element values (`asort`) or the index values (`asorti`). Here's an example of using `asort`:

```
$ gawk 'BEGIN{
> var["a"] = 1
> var["g"] = 2
> var["m"] = 3
> var["u"] = 4
> asort(var, test)
> for (i in test)
>     print "Index:",i," - value:",test[i]
> }'
Index: 4 - value: 4
Index: 1 - value: 1
Index: 2 - value: 2
Index: 3 - value: 3
$
```

The new array, `test`, contains the newly sorted data elements of the original array, but the index values are now changed to numerical values, indicating the proper sort order.

The `split` function is a great way to push data fields into an array for further processing:

```
$ gawk 'BEGIN{ FS="," }{
> split($0, var)
> print var[1], var[5]
> }' data1
data11 data15
data21 data25
data31 data35
$
```

The new array uses sequential numbers for the array index, starting with index value 1 containing the first data field.

Time Functions

The gawk programming language contains a few functions to help you deal with time values, shown in [Table 21.6](#).

Table 21.6 The gawk Time Functions

Function	Description
<code>mktime (datespec)</code>	Converts a date specified in the format YYYYMMDD HH MMSS [DST] into a timestamp value.
<code>strftime (format [, timestamp])</code>	Formats either the current time of day timestamp, or timestamp if provided, into a formatted day and date, using the <code>date()</code> shell function format.
<code>systemtime()</code>	Returns the timestamp for the current time of day.

The time functions are often used when working with log files that contain dates that you need to compare. By converting the text representation of a date to the epoch time (the number of seconds since midnight, January 1, 1970) you can easily compare dates.

The following is an example of using the time functions in a gawk program:

```
$ gawk 'BEGIN{
> date = systemtime()
> day = strftime("%A, %B %d, %Y", date)
> print day
> }'
Friday, December 28, 2010
$
```

This example uses the `systemtime` function to retrieve the current epoch timestamp from the system and then uses the `strftime` function to convert it into a human-readable format using the `date` shell command's `date` format

characters.

User-Defined Functions

You're not limited to just using the built-in functions available in gawk. You can create your own functions for use in your gawk programs. This section shows you how to define and use your own functions in your gawk programs.

Defining a Function

To define your own function, you must use the `function` keyword:

```
function name([variables])  
{  
    statements  
}
```

The function name must uniquely identify your function. You can pass one or more variables into the function from the calling gawk program:

```
function printthird()  
{  
    print $3  
}
```

This function will print the third data field in the record.

The function can also return a value using the `return` statement:

```
return value
```

The value can be a variable, or an equation that evaluates to a value:

```
function myrand(limit)
{
    return int(limit * rand())
}
```

You can assign the value returned from the function to a variable in the gawk program:

```
x = myrand(100)
```

The variable will contain the value returned from the function.

Using Your Functions

When you define a function, it must appear by itself before you define any programming sections (including the `BEGIN` section). This may look a little odd at first, but it helps keep the function code separate from the rest of the gawk program:

```
$ gawk `
> function myprint()
> {
>     printf "%-16s - %s\n", $1, $4
> }
> BEGIN{FS="\n"; RS=""}
> {
>     myprint()
> }' data2
Riley Mullen      - (312)555-1234
Frank Williams    - (317)555-9876
Haley Snell       - (313)555-4938
$
```

The function defines the `myprint()` function, which formats the first and fourth data fields in the record for printing. The gawk program then uses the function to display the data from the data file.

Once you define a function, you can use it as often as

necessary in the program section of the code. This saves lots of work when using long algorithms.

Creating a Function Library

Obviously, having to rewrite your gawk functions every time you need them is not all that pleasant of an experience. However, gawk provides a way for you to combine your functions into a single library file that you can use in all of your gawk programming.

First, you need to create a file that contains all of your gawk functions:

```
$ cat funclib
function myprint()
{
    printf "%-16s - %s\n", $1, $4
}
function myrand(limit)
{
    return int(limit * rand())
}
function printthird()
{
    print $3
}
$
```

The `funclib` file contains three function definitions. To use them, you need to use the `-f` command line parameter. Unfortunately, you can't combine the `-f` command line parameter with an inline gawk script, but you can use multiple `-f` parameters on the same command line.

Thus, to use your library, just create a file that contains your gawk program, and specify both the library file and your program file on the command line:

```
$ cat script4
BEGIN{ FS="\n"; RS="" }
{
    myprint()
}
$ gawk -f funclib -f script4 data2
Riley Mullen      - (312)555-1234
Frank Williams    - (317)555-9876
Haley Snell       - (313)555-4938
$
```

Now all you need to do is add the `funclib` file to your gawk command line whenever you need to use a function defined in the library.

Summary

This chapter walked you through the more advanced features of the gawk programming language. Every programming language requires using variables, and gawk is no different. The gawk programming language includes some built-in variables that you can use to reference specific data field values and retrieve information about the number of data fields and records processed in the data file. You can also create your own variables for use in your scripts.

The gawk programming language also provides many of the standard structured commands you'd expect from a programming language. You can easily create fancy programs using `if-then` logic, `while`, and `do-while` and `for` loops. Each of these commands allows you to alter the flow of your gawk program script to iterate through data field

values to create detailed data reports.

The `printf` command is a great tool to have if you need to customize your report output. It allows you to specify the exact format for displaying data from the gawk program script. You can easily create formatted reports, placing data elements in exactly the correct position.

Finally, this chapter discussed the many built-in functions available in the gawk programming language and showed you how to create your own functions. The gawk program contains many useful functions for handling mathematical features, such as standard square roots and logarithms, as well as trigonometric functions. There are also several string-related functions that make extracting substrings from larger strings a breeze.

You aren't limited to the built-in functions in the gawk program. If you're working on an application that uses lots of specialized algorithms, you can create your own functions to process the algorithms, and then use those functions in your own code. You can also set up a library file containing all of the functions you use in your gawk programs, saving you time and effort in all of your coding.

The next chapter switches gears a little. It examines a few other shell environments you may run into in your Linux shell-scripting endeavors. While the bash shell is the most common shell used in Linux, it's not the only shell. It helps to know a little about some of the other shells available and how they differ from the bash shell.

Chapter 22

Working with Alternative Shells

In This Chapter

- Understanding the dash shell
- Programming in the dash shell
- Introducing the zsh shell
- Writing scripts for zsh

While the bash shell is the most widely used shell in Linux distributions, it's not the only one. Now that you've seen the standard Linux bash shell, and what you can do with it, it's time to examine a few other shells available in the Linux world. This chapter describes two other shells that you may run into in your Linux journey and how they differ from the bash shell.

What Is the dash Shell?

The Debian dash shell has had an interesting past. It's a direct descendant of the ash shell, a simple copy of the original Bourne shell available on Unix systems (see

Chapter 1). Kenneth Almquist created a small-scale version of the Bourne shell for Unix systems and called it the Almquist shell, which was then shortened to *ash*. This original version of the ash shell was extremely small and fast but without many advanced features, such as command line editing or history features, making it difficult to use as an interactive shell.

The NetBSD Unix operating system adopted the ash shell and still uses it today as the default shell. The NetBSD developers customized the ash shell by adding several new features, making it closer to the Bourne shell. The new features include command line editing using both emacs and vi editor commands, as well as a history command to recall previously entered commands. This version of the ash shell is also used by the FreeBSD operating system as the default login shell.

The Debian Linux distribution created its own version of the ash shell (called Debian ash, or *dash*) for inclusion in its version of Linux. For the most part, dash copies the features of the NetBSD version of the ash shell, providing the advanced command line editing capabilities.

However, to add to the shell confusion, the dash shell is actually not the default shell in many Debian-based Linux distributions. Because of the popularity of the bash shell in Linux, most Debian-based Linux distributions use the bash shell as the normal login shell, and only use the dash shell as a quick-start shell for the installation script to install the distribution files.

The exception is the popular Ubuntu distribution. This often confuses shell script programmers and causes a great number of problems with running shell scripts in a Linux environment. The Ubuntu Linux distribution uses the bash shell as the default interactive shell, but uses the dash shell as the default `/bin/sh` shell. This “feature” really confuses shell script programmers.

As you saw in Chapter 10, every shell script must start with a line that declares the shell used for the script. In our bash shell scripts, we've been using the following:

```
#!/bin/bash
```

This tells the shell to use the shell program located at `/bin/bash` to execute the script. In the Unix world, the default shell was always `/bin/sh`. Many shell script programmers familiar with the Unix environment copy this into their Linux shell scripts:

```
#!/bin/sh
```

On most Linux distributions, the `/bin/sh` file is a symbolic link (see Chapter 3) to the `/bin/bash` shell program. This allows you to easily port shell scripts designed for the Unix Bourne shell to the Linux environment without having to modify them.

Unfortunately, the Ubuntu Linux distribution links the `/bin/sh` file to the `/bin/dash` shell program. Because the dash shell contains only a subset of the commands available in the original Bourne shell, this can (and often does) cause some shell scripts to not work properly.

The next section walks you through the basics of the

dash shell and how it differs from the bash shell. This is especially important to know if you write bash shell scripts that may need to be run in an Ubuntu environment.

The dash Shell Features

While both the bash shell and the dash shell are modeled after the Bourne shell, they have some differences. This section walks you through the features found in the Debian dash shell to acquaint you with how the dash shell works before we dive into the shell scripting features.

The dash Command Line Parameters

The dash shell uses command line parameters to control its behavior. [Table 22.1](#) lists the command line parameters and describes what each one does.

Table 22.1 The dash Command Line Parameters

Parameter	Description
-a	Export all variables assigned to the shell.
-c	Read commands from a specified command string.
-e	If not interactive, exit immediately if any untested command fails.
-f	Display pathname wildcard characters.
-n	If not interactive, read commands but don't execute them.
-u	Write an error message to STDERR when attempting to expand a variable that is not set.

-v	Write input to STDERR as it is read.
-x	Write each command to STDERR as it's executed.
-I	Ignore EOF characters from the input when in interactive mode.
-i	Force the shell to operate in interactive mode.
-m	Turn on job control (enabled by default in interactive mode).
-s	Read commands from STDIN (the default behavior if no file arguments are present).
-E	Enable the emacs command line editor.
-v	Enable the vi command line editor.

There are just a few additional command line parameters that Debian added to the original ash shell command line parameter list. The `-E` and `-v` command line parameters enable the special command line editing features of the dash shell.

The `-E` command line parameter allows you to use the emacs editor commands for editing command line text (see Chapter 9). You can use all of the emacs commands for manipulating text on a single line using the Ctrl and Meta key combinations.

The `-v` command line parameter allows you to use the vi editor commands for editing command line text (again, see Chapter 9). This feature allows you to switch between normal mode and vi editor mode on the command line by using the `ESC` key. When you're in vi editor mode, you can use all of the standard vi editor commands (such as `x` to delete a character, and `i` to insert text). Once you are finished editing the command line, you must hit the Esc key

again to exit vi editor mode.

The dash Environment Variables

There are quite a few default environment variables that the dash shell uses to track information, and you can create your own environment variables as well. This section describes the environment variables and how dash handles them.

Default Environment Variables

[Table 22.2](#) shows the default dash environment variables and describes what they are used for.

[Table 22.2](#) The dash Shell Environment Variables

Variable	Description
CDPATH	The search path used for the <code>cd</code> command.
HISTSIZE	The number of lines stored in the history file.
HOME	The user's default login directory.
IFS	The input field separator characters. The default value is a space, tab, and newline.
MAIL	The name of the user's mailbox file.
MAILCHECK	The frequency with which to check the mailbox file for new mail.
MAILPATH	A colon-separated list of multiple mailbox filenames. If set, this value overrides the <code>MAIL</code> environment variable.
OLDPWD	The value of the previous working directory.
PATH	The default search path for executable files.
PPID	The Process ID of the current shell's parent.

PS1	The primary shell command line interface prompt.
PS2	The secondary shell command line interface prompt.
PS4	A character printed before each line when an execution trace is enabled.
PWD	The value of the current working directory.
TERM	The default terminal setting for the shell.

You should notice that the dash environment variables are very similar to the environment variables used in bash (see Chapter 5). This is not by accident. Remember that both the dash and bash shells are extensions of the Bourne shell, so they both incorporate many of its features. However, because of its goal of simplicity, the dash shell contains significantly fewer environment variables than the bash shell. You need to take this into consideration when creating shell scripts in an ash shell environment.

The dash shell uses the `set` command to display environment variables:

```
$set
COLORTERM=''
DESKTOP_SESSION='default'
DISPLAY=':0.0'
DM_CONTROL='/var/run/xdmctl'
GS_LIB='/home/atest/.fonts'
HOME='/home/atest'
IFS='
'
KDEROOTHOME='/root/.kde'
KDE_FULL_SESSION='true'
KDE_MULTIHEAD='false'
KONSOLE_DCOP='DCOPRef(konsole-5293,konsole)'
KONSOLE_DCOP_SESSION='DCOPRef(konsole-5293,session-1)'
LANG='en_US'
LANGUAGE='en'
LC_ALL='en_US'
LOGNAME='atest'
OPTIND='1'
```

```
PATH=/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin
PPID='5293'
PS1='$ '
PS2='> '
PS4='+ '
PWD='/home/atest'
SESSION_MANAGER='local/testbox:/tmp/.ICE-unix/5051'
SHELL='/bin/dash'
SHLVL='1'
TERM='xterm'
USER='atest'
XCURSOR_THEME='default'
_='ash'
$
```

Your default dash shell environment will most likely differ, as different Linux distributions assign different default environment variables at login.

Positional Parameters

In addition to the default environment variables, the dash shell also assigns special variables to any parameters defined in the command line. Here are the positional parameter variables available for use in the dash shell:

- `$0`: The name of the shell
- `$n`: The *n*th position parameter
- `$*`: A single value with the contents of all of the parameters, separated by the first character in the IFS environment variable, or a space if IFS isn't defined
- `$@`: Expands to multiple arguments consisting of all the command line parameters
- `$#`: The number of positional parameters
- `$?`: The exit status of the most recent command

- `$-:` The current option flags
- `$$:` The process ID (PID) of the current shell
- `$!:` The process ID (PID) of the most recent background command

All of the dash positional parameters mimic the same positional parameters available in the bash shell. You can use each of the positional parameters in your shell scripts just as you would in the bash shell.

User-Defined Environment Variables

The dash shell also allows you to set your own environment variables. As with bash, you can define a new environment variable on the command line by using the `assign` statement:

```
$ testing=10
$ echo $testing
10
$
```

By default, environment variables are only visible in the shell session in which they're defined. To allow an environment variable to be visible in a child shell or process, you must use the `export` command:

```
$ testing=10 ; export testing
$ dash
$ echo $testing
10
$
```

Without the `export` command, user-defined environment variables are only visible in the current shell or process.

Caution

There's one huge difference between dash variables and bash variables. The dash shell doesn't support variable arrays. This small feature causes all sorts of problems for advanced shell script writers.

The dash Built-in Commands

Just as with the bash shell, the dash shell contains a set of built-in commands that it recognizes. You can use these commands directly from the command line interface or you can incorporate them in your shell scripts. [Table 22.3](#) lists the dash shell built-in commands.

Table 22.3 The dash Shell Built-In Commands

Command	Description
alias	Create an alias string to represent a text string.
bg	Continue specified job in background mode.
cd	Switch to the specified directory.
echo	Display a text string and environment variables.
eval	Concatenate all arguments with a space.
exec	Replace the shell process with the specified command.
exit	Terminate the shell process.
export	Export the specified environment variable for use in all child shells.
fc	List, edit, or re-execute commands previously entered on the command line.
fg	Continue specified job in foreground mode.

getopts	Obtain options and arguments from a list of parameters.
hash	Maintain and retrieve a hash table of recent commands and their locations.
pwd	Display the value of the current working directory.
read	Read a line from STDIN and assign the value to a variable.
readonly	Read a line from STDIN to a variable that can't be changed.
printf	Display text and variables using a formatted string.
set	List or set option flags and environment variables.
shift	Shift the positional parameters a specified number of times.
test	Evaluate an expression and return 0 if true, or 1 if false.
times	Display the accumulated user and system times for the shell and all shell processes.
trap	Parse and execute an action when the shell receives a specified signal.
type	Interpret the specified name and display the resolution (alias, built-in, command, keyword).
ulimit	Query or set limits on processes.
umask	Set the value of the default file and directory permissions.
unalias	Remove the specified alias.
unset	Remove the specified variable or option flag from the exported variables.
wait	Wait for the specified job to complete and return the exit status.

You probably recognize all of these built-in commands from the bash shell. The dash shell supports many of the same built-in commands as the bash shell. You'll notice that there aren't any commands for the command history file, nor for the directory stack. The dash shell doesn't support

these features.

Scripting in dash

Unfortunately, the dash shell doesn't recognize all of the scripting features of the bash shell. Shell scripts written for the bash environment often fail when run in the dash shell, causing all sorts of grief for shell script programmers. This section describes the differences you'll need to be aware of to get your shell scripts to run properly in a dash shell environment.

Creating dash Scripts

You probably guessed by now that creating shell scripts for the dash shell is pretty similar to creating shell scripts for the bash shell. You should always specify which shell you want to use in your script to ensure that the script runs with the proper shell.

You do this on the first line of the shell:

```
#!/bin/dash
```

You can also specify a shell command line parameter on this line, as was demonstrated earlier in “The dash Command Line Parameters” section.

Things That Won't Work

Unfortunately, because the dash shell is only a subset of the

Bourne shell features, there are a few things in bash shell scripts that won't work in the dash shell. These are often called *bashisms*. This section is a quick summary of bash shell features you may be used to using in your bash shell scripts that won't work if you're in a dash shell environment.

Using Arithmetic

Chapter 10 showed three ways to express a mathematical operation in the bash shell script:

- **Using the `expr` command:** `expr operation`
- **Using square brackets:** `$(operation)`
- **Using double parentheses:** `$((operation))`

The dash shell supports the `expr` command and the double parentheses method but doesn't support the square bracket method. This can be a problem if you've got lots of mathematical operations that use the square brackets:

```
$ cat test5
#!/bin/dash
# testing mathematical operations

value1=10
value2=15

value3=$( $value1 * $value2 ]
echo "The answer is $value3"
$ ./test5
./test5: 7: value1: not found
The answer is
$
```

The proper format for performing mathematical operations in dash shell scripts is to use the double parentheses method:

```
$ cat test5b
```

```
#!/bin/dash
# testing mathematical operations

value1=10
value2=15

value3=$(( $value1 * $value2 ))
echo "The answer is $value3"
$ ./test5b
The answer is 150
$
```

Now the shell can perform the calculation properly.

The test Command

While the dash shell supports the `test` command, you must be careful how you use it. The bash shell version of the `test` command is slightly different than the dash shell version.

The bash shell `test` command allows you to use the double equal sign (`==`) to test if two strings are equal. This is an add-on to accommodate programmers familiar with using this format in other programming languages:

```
$ cat test6
#!/bin/bash
# testing the == comparison

test1=abcdef
test2=abcdef

if [ $test1 == $test2 ]
then
    echo "They're the same!"
else
    echo "They're different"
fi
$ ./test6
They're the same!
$
```

Simple enough. However, if you run this script in a dash

shell environment, you'll get an unwelcome outcome:

```
$ ./test6
[: ==: unexpected operator
They're different
$
```

The `test` command available in the `dash` shell doesn't recognize the `==` symbol for text comparisons. Instead, it only recognizes the `=` symbol. If you change the text comparison symbol to just a single equal sign, things are just fine in both the `bash` and `dash` shell environments:

```
$ cat test7
#!/bin/dash
# testing the = comparison

test1=abcdef
test2=abcdef

if [ $test1 = $test2 ]
then
    echo "They're the same!"
else
    echo "They're different"
fi
$ ./test7
They're the same!
$
```

This little bashism is responsible for many hours of frustration for shell programmers!

The echo Statement Options

The simple `echo` statement is also another source of frustration for `dash` shell programmers. It doesn't behave the same way in the `dash` and `bash` shells.

In the `bash` shell, if you want to display a special character in the output, you must use the `-e` command line

parameter:

```
echo -e "This line contains\t a special character"
```

The `echo` statement contains the `\t` special character to represent a tab. Without the `-e` command line parameter, the `bash` version of the `echo` statement ignores the special character. Here's a test that demonstrates this in the `bash` shell:

```
$ cat test8
#!/bin/bash
# testing echo commands

echo "This is a normal test"
echo "This test uses a\t special character"
echo -e "This test uses a\t special character"
echo -n "Does this work: "
read test
echo "This is the end of the test"
$ ./test8
This is a normal test
This test uses a\t special character
This test uses a special character
Does this work: N
This is the end of the test
$
```

The `echo` statement that didn't use the `-e` command line parameter just displayed the `\t` character as normal text, so to get the tab character, you must use the `-e` command line parameter.

In the `dash` shell, things are a little different. The `echo` statement in the `dash` shell automatically recognizes and displays special characters. Because of this, there is no `-e` command line parameter. If you try to run the same script in a `dash` environment, you get this output:

```
$ ./test8
This is a simple test
This line uses a      special character
```

```
e This line uses a      special character
Does this work: N
This is the end of the test
$
```

As you can see from the output, the dash shell version of the `echo` command recognizes the special character in the line without the `-e` command line parameter just fine, but for the line that does contain the `-e` command line parameter, it gets confused and displays the `-e` as normal text.

Unfortunately, there's no simple solution to this problem. If you must write scripts that work in both the bash and dash shell environments, the best solution is to use the `printf` command to display text. This command works the same way in both shell environments, and it can display special characters just fine.

The function Command

Chapter 16 showed you how to define your own functions in your shell scripts. The bash shell supports two methods for defining functions. The first method is to use the `function` statement:

```
function name {
    commands
}
```

The `name` attribute defines a unique name assigned to the function. Each function you define in your script must be assigned a unique name.

The `commands` are one or more bash shell commands that make up your function. When you call the function, the bash shell executes each of the commands in the order they

appear in the function, just as in a normal script.

The second format for defining a function in a bash shell script more closely follows how functions are defined in other programming languages:

```
name() {  
  commands  
}
```

The empty parentheses after the function name indicate that you're defining a function. The same naming rules apply in this format as the original shell script function format.

The dash shell doesn't support the first method of defining functions (it doesn't support the `function` statement). Instead, in the dash shell you must define a function using the function name with parentheses. If you try to run a function designed for the bash shell in a dash shell, you'll get an error message:

```
$ cat test9  
#!/bin/dash  
# testing functions  
  
function func1() {  
    echo "This is an example of a function"  
}  
count=1  
while [ $count -le 5 ]  
do  
    func1  
    count=$(( $count + 1 ))  
done  
echo "This is the end of the loop"  
func1  
echo "This is the end of the script"  
$ ./test9  
./test9: 4: Syntax error: "(" unexpected  
$
```


Instead of assigning the function code to the function, the dash shell executed the code within the function definition and then complained about the format of the shell script.

If you're writing shell scripts that may be used in the dash environment, always use the second method of defining your functions:

```
$ cat test10
#!/bin/dash
# testing functions

func1() {
    echo "This is an example of a function"
}

count=1
while [ $count -le 5 ]
do
    func1
    count=$(( $count + 1 ))
done
echo "This is the end of the loop"
func1
echo "This is the end of the script"
$ ./test10
This is an example of a function
This is an example of a function
This is an example of a function
This is an example of a function
This is an example of a function
This is the end of the loop
This is an example of a function
This is the end of the script
$
```

Now the dash shell recognized the function defined in the script just fine and was able to use it within the script.

The zsh Shell

Another popular shell that you may run into is the Z shell (called zsh). The zsh shell is an open source Unix shell developed by Paul Falstad. It takes ideas from all of the existing shells and adds many unique features to create a full-blown advanced shell designed for programmers.

The following are some of the features that make the zsh shell unique:

- Improved shell option handling
- Shell compatibility modes
- Loadable modules

Of all these features, a loadable module is the most advanced feature in shell design. As you've seen in the bash and dash shells, each shell contains a set of built-in commands that are available without the need for external utility programs. The benefit of built-in commands is execution speed. The shell doesn't have to load a utility program into memory before running it; the built-in commands are already in the shell memory, ready to go.

The zsh shell provides a core set of built-in commands, plus the capability to add additional *command modules*. Each command module provides a set of additional built-in commands for specific circumstances, such as network support and advanced math functions. You can add only the modules you think you need for your specific situation.

This feature provides a great way to either limit the size of the zsh shell for situations that require a small shell size and few commands or expand the number of available built-in commands for situations that require faster

execution speeds.

Parts of the zsh Shell

This section walks you through the basics of the zsh shell, showing the built-in commands that are available (or can be added by installing modules), as well as the command line parameters and environment variables used by the zsh shell.

Shell Options

Most shells use command line parameters to define the behavior of the shell. The zsh shell uses a few command line parameters to define the operation of the shell, but mostly it uses *options* to customize the behavior of the shell. You can set shell options either on the command line or within the shell itself using the `set` command.

[Table 22.4](#) lists the command line parameters available for the zsh shell.

[Table 22.4](#) The zsh Shell Command Line Parameters

Parameter	Description
-c	Execute only the specified command and exit.
-i	Start as an interactive shell, providing a command line interface prompt.
-s	Force the shell to read commands from <code>STDIN</code> .
-o	Specify command line options.

While this may seem like a small set of command line parameters, the `-o` parameter is somewhat misleading. It allows you to set shell options that define features within the shell. By far the `zsh` shell is the most customizable shell available. There are lots of features that you can alter for your shell environment. The different options fit into several general categories:

- **Changing directories:** Options that control how the `cd` and `dirs` commands handle directory changes
- **Completion:** Options that control command completion features
- **Expansion and globbing:** Options that control file expansion in commands
- **History:** Options that control command history recall
- **Initialization:** Options that control how the shell handles variables and startup files when started
- **Input/Output:** Options that control command handling
- **Job Control:** Options that dictate how the shell handles and starts jobs
- **Prompting:** Options that define how the shell works with command line prompts
- **Scripts and functions:** Options that control how the shell processes shell scripts and defines shell functions
- **Shell emulation:** Options that allow you to set the behavior of the `zsh` shell to mimic the behavior of

other shell types

- **Shell state:** Options that define what type of shell to start
- **zle:** Options for controlling the zsh line editor (zle) feature
- **Option aliases:** Special options that can be used as aliases for other option names

With this many different categories of shell options, you can imagine just how many actual options the zsh shell supports. The following sections show a sampling of the different zsh shell options available for you to use when customizing your zsh shell environment.

Shell State

There are six different zsh shell options that define the type of shell to start:

- **interactive** (`-i`): Provides a command line interface prompt for entering built-in commands and program names.
- **login** (`-l`): The default zsh shell type, processes the zsh shell startup files, and provides a command line interface prompt.
- **privileged** (`-p`): The default if the effective user ID (UID) of the user is not the same as the real UID (the user has become the root user). This option disables the user startup files.
- **restricted** (`-r`): Restricts the user to a specified

directory structure in the shell.

- **shin_stdin** (-s): Commands are read from STDIN.
- **single_command** (-t): Executes a single command from STDIN and exits.

The shell states define whether or not the shell starts with a command line interface prompt and what access the user has within the shell.

Shell Emulation

The shell emulation options allow you to customize the zsh shell to perform similar to the C shell (csh) or Korn shell (ksh), two other shells popular with programmers. These options are as follows:

- **bsd_echo**: Make the `echo` statement compatible with the C shell `echo` command.
- **csh_junkie_history**: Use the `history` command without a specifier to reference the previous command.
- **csh_junkie_loops**: Allow `while` and `for` loops to use `end` like the C shell instead of `do` and `done`.
- **csh_junkie_quotes**: Change the rules of using single and double quotation marks to match the C shell.
- **csh_nullcmd**: Don't use the values of the `NULLCMD` and `READNULLCMD` variables when executing redirections with no commands.
- **ksh_arrays**: Use Korn-style arrays, starting numeric

arrays at 0, and require braces to reference array elements.

- **ksh_autoload:** Emulate the Korn shell autoload function feature.
- **ksh_option_print:** Emulate the Korn shell method of printing options.
- **ksh_typeset:** Alter the way that the typeset command arguments are processed.
- **posix_builtins:** Use the `builtin` command to execute built-in commands.
- **sh_file_expansion:** Perform file name expansion before any other expansion.
- **sh_nullcmd:** Don't use the `NULLCMD` and `READNULLCMD` variables when performing redirections.
- **sh_option_letters:** Interpret single letter shell command line options similar to the Korn shell.
- **sh_word_split:** Perform field splitting on unquoted parameter expansions.
- **traps_async:** While waiting for a program to exit, handle signals and run traps immediately.

By having multiple options you can pick and choose which csh or ksh shell feature you need to emulate in your zsh shell, instead of having to emulate the entire shell.

Initialization

There are a few options for handling the shell startup features:

- **all_export:** All parameters and variables are exported to child shell processes automatically.
- **global_export:** Parameters exported to the environment will not be made local to the function.
- **global_rcs:** If not set, the zsh shell doesn't run the global startup files, but it still runs local startup files.
- **rcs:** If not set, the zsh shell runs the `/etc/zshenv` startup file, but no others.

The initialization options allow you to specify which (if any) zsh shell startup files are run in your shell environment. You can also set these values within startup files themselves to limit which ones the shell executes.

Scripts and Functions

The script and function options allow you to customize the shell scripting environment in the zsh shell. This is a handy way to set the way functions perform within the shell.

- **c_bases:** Display hexadecimal numbers in C format (0xdddd) instead of shell format (16#dddd).
- **err_exit:** If a command exits with a non-zero exit status, perform the command in the ZERR trap and exit.
- **err_return:** If a command has a non-zero exit status, return immediately from the enclosing function.
- **eval_lineno:** If set, the line numbers of expressions evaluated using the `eval` built-in command are tracked separately from the rest of the shell environment.

- **exec:** Execute commands. If this option isn't set, read commands and report errors, but don't execute the commands.
- **function_argzero:** Set `$0` to the name of the function or script.
- **local_options:** If set, when a shell function returns, all of the options that were set before the function are restored.
- **local_traps:** If set, when a signal trap is set within a function, the previous status of the trap is restored when the function exits.
- **multios:** Perform implicit `tee` or `cat` commands when multiple redirections are attempted.
- **octal_zeroes:** Interpret any integer string starting with a zero as an octal number.
- **typeset_silent:** If not set, using `typeset` with the name of a parameter displays the current value of the parameter.
- **verbose:** Displays shell input lines as they are read by the shell.
- **xtrace:** Displays commands and their arguments as the shell executes them.

The zsh shell allows you to customize lots of features that occur when you're exiting functions defined in the shell.

Built-in Commands

The zsh shell is unique in that it allows you to expand the built-in commands available in the shell. This provides for a wealth of speedy utilities at your fingertips for a host of different applications.

This section describes the core built-in commands, along with the various modules available at the time of this writing.

Core Built-in Commands

The core of the zsh shell contains the basic built-in commands you're used to seeing in other shells. [Table 22.5](#) describes the built-in commands available for you.

Table 22.5 The zsh Core Built-In Commands

Command	Description
<code>alias</code>	Define an alternate name for a command and arguments.
<code>autoload</code>	Preload a shell function into memory for quicker access.
<code>bg</code>	Execute a job in background mode.
<code>bindkey</code>	Bind keyboard combinations to commands.
<code>builtin</code>	Execute the specified built-in command instead of an executable file of the same name.
<code>bye</code>	The same as <code>exit</code> .
<code>cd</code>	Change the current working directory.
<code>chdir</code>	Change the current working directory.
<code>command</code>	Execute the specified command as an external file instead of a function or built-in command.
<code>declare</code>	Set the data type of a variable (same as <code>typeset</code>).
<code>dirs</code>	Display the contents of the directory stack.

disable	Temporarily disable the specified hash table elements.
disown	Remove the specified job from the job table.
echo	Display variables and text.
emulate	Set zsh to emulate another shell, such as the Bourne, Korn, or C shells.
enable	Enable the specified hash table elements.
eval	Execute the specified command and arguments in the current shell process.
exec	Execute the specified command and arguments replacing the current shell process.
exit	Exit the shell with the specified exit status. If none specified, use the exit status of the last command.
export	Allow the specified environment variable names and values to be used in child shell processes.
false	Return an exit status of 1.
fc	Select a range of commands from the history list.
fg	Execute the specified job in foreground mode.
float	Set the specified variable for use as a floating point variable.
functions	Set the specified name as a function.
getln	Read the next value in the buffer stack and place it in the specified variable.
getopts	Retrieve the next valid option in the command line arguments and place it in the specified variable.
hash	Directly modify the contents of the command hash table.
history	List the commands contained in the history file.
integer	Set the specified variable for use as an integer value.
jobs	List information about the specified job, or all jobs assigned to the shell process.

kill	Send a signal (Default <code>SIGTERM</code>) to the specified process or job.
let	Evaluate a mathematical operation and assign the result to a variable.
limit	Set or display resource limits.
local	Set the data features for the specified variable.
log	Display all users currently logged in who are affected by the watch parameter.
logout	Same as <code>exit</code> , but only works when the shell is a login shell.
popd	Remove the next entry from the directory stack.
print	Display variables and text.
printf	Display variables and text using C-style format strings.
pushd	Change the current working directory, and put the previous directory in the directory stack.
pushln	Place the specified arguments into the editing buffer stack.
pwd	Display the full pathname of the current working directory.
read	Read a line and assign data fields to the specified variables using the <code>IFS</code> characters.
readonly	Assign a value to a variable that can't be changed.
rehash	Rebuild the command hash table.
set	Set options or positional parameters for the shell.
setopt	Set the options for a shell.
shift	Read and delete the first positional parameter, and then shift the remaining ones down one position.
source	Find the specified file and copy its contents into the current location.
suspend	Suspend the execution of the shell until it receives a <code>SIGCONT</code> signal.
test	Return an exit status of 0 if the specified condition is <code>TRUE</code> .

times	Display the cumulative user and system times for the shell and processes that run in the shell.
trap	Block the specified signals from being processed by the shell, and execute the specified commands if the signals are received.
true	Return a zero exit status.
ttctl	Lock and unlock the display.
type	Display how the specified command would be interpreted by the shell.
typeset	Set or display attributes of variables.
ulimit	Set or display resource limits of the shell or processes running in the shell.
umask	Set or display the default permissions for creating files and directories.
unalias	Remove the specified command alias.
unfunction	Remove the specified defined function.
unhash	Remove the specified command from the hash table.
unlimit	Remove the specified resource limit.
unset	Remove the specified variable attribute.
unsetopt	Remove the specified shell option.
wait	Wait for the specified job or process to complete.
whence	Display how the specified command would be interpreted by the shell.
where	Display the pathname of the specified command if found by the shell.
which	Display the pathname of the specified command using csh-style output.
zcompile	Compile the specified function or script for faster autoloading.
zmodload	Perform operations on loadable zsh modules.

The zsh shell is no slouch when it comes to providing built-in commands! You should recognize most of these commands from their bash counterparts. The most important features of the zsh shell built-in commands are modules.

Add-in Modules

There's a long list of modules that provide additional built-in commands for the zsh shell, and the list continues to grow as resourceful programmers create new modules. [Table 22.6](#) shows the modules available at the time of this writing.

[Table 22.6](#) The zsh Modules

Module	Description
zsh/cap	POSIX compatibility commands
zsh/clone	Commands to clone a running shell to another terminal
zsh/compctl	Commands to control command completion
zsh/complete	Command line completion commands
zsh/complist	Commands for command line completion listing extensions
zsh/computil	Utility commands for command line completion
zsh/datetime	Additional date and time commands and variables
zsh/deltochar	A line editor function replicating emacs functionality
zsh/files	Commands for basic file handling
zsh/mapfile	Access to external files via associative arrays
zsh/mathfunc	Additional scientific functions
zsh/parameter	Access to command hash tables via associative arrays

<code>zsh/pcore</code>	The extended regular expression library
<code>zsh/sched</code>	Scheduling commands for providing timed command execution
<code>zsh/net/socket</code>	Unix domain socket support
<code>zsh/stat</code>	Access to the <code>stat</code> system call to provide system statistics
<code>zsh/system</code>	Interface for various low-level system features
<code>zsh/net/tcp</code>	Access to TCP sockets
<code>zsh/termcap</code>	Interface to the <code>termcap</code> database
<code>zsh/terminfo</code>	Interface to the <code>terminfo</code> database
<code>zsh/zftp</code>	A specialized FTP client command
<code>zsh/zle</code>	The <code>zshell</code> line editor
<code>zsh/zleparameter</code>	Access to modify <code>zle</code> using variables
<code>zsh/zprof</code>	Allow profiling for shell functions
<code>zsh/zpty</code>	Start a command in a pseudo-terminal
<code>zsh/zselect</code>	Block and return when file descriptors are ready
<code>zsh/zutil</code>	Various shell utilities

The `zsh` shell modules cover a wide range of topics, from providing simple command line editing features to advanced networking functions. The idea behind the `zsh` shell is to provide a basic minimum shell environment and let you add on the pieces you need to accomplish your programming job.

Viewing, Adding, and Removing Modules

The `zmodload` command is the interface to the `zsh` modules.

You use this command to view, add, and remove modules from the zsh shell session.

Using the `zmodload` command without any command line parameters displays the currently installed modules in your zsh shell:

```
% zmodload
zsh/zutil
zsh/complete
zsh/main
zsh/terminfo
zsh/zle
zsh/parameter
%
```

Different zsh shell implementations include different modules by default. To add a new module, just specify the module name on the `zmodload` command line:

```
% zmodload zsh/zftp
%
```

Nothing indicates that the module loaded. You can perform another `zmodload` command, and the new module should appear in the list of installed modules.

Once you load a module, the commands associated with the module are available as built-in commands:

```
% zftp open myhost.com rich testing1
Welcome to the myhost FTP server.
% zftp cd test
% zftp dir
01-21-11 11:21PM      120823 test1
01-21-11 11:23PM      118432 test2
% zftp get test1 > test1.txt
% zftp close
%
```

The `zftp` command allows you to conduct a complete FTP session directly from your zsh shell command line! You can

incorporate these commands into your zsh shell scripts to perform file transfers directly from your scripts.

To remove an installed module, use the `-u` parameter, along with the module name:

```
% zmodload -u zsh/zftp
% zftp
zsh: command not found: zftp
%
```

Note

It's a common practice to place `zmodload` commands in the `$HOME/.zshrc` startup file so that your favorite functions load automatically when the zsh shell starts.

Scripting with zsh

The main purpose of the zsh shell was to provide an advanced programming environment for shell programmers. With that in mind, it's no surprise that the zsh shell offers many features that make shell scripting easier.

Mathematical Operations

As you would expect, the zsh shell allows you to perform mathematical functions with ease. In the past, the Korn shell has led the way in supporting mathematical operations by providing support for floating-point numbers. The zsh shell has full support for floating-point numbers in all of its

mathematical operations!

Performing Calculations

The zsh shell supports two methods for performing mathematical operations:

- The `let` command
- Double parentheses

When you use the `let` command, you should enclose the operation in double quotation marks to allow for spaces:

```
% let value1=" 4 * 5.1 / 3.2 "  
% echo $value1  
6.374999999999999991  
%
```

Notice that using floating point numbers introduces a precision problem. To solve this, it's always a good idea to use the `printf` command, and specify the decimal precision needed to correctly display the answer:

```
% printf "%6.3f\n" $value1  
6.375  
%
```

Now that's much better!

The second method is to use the double parentheses. This method incorporates two techniques for defining the mathematical operation:

```
% value1=$(( 4 * 5.1 ))  
% (( value2 = 4 * 5.1 ))  
% printf "%6.3f\n" $value1 $value2  
20.400  
20.400  
%
```

Notice that you can place the double parentheses either

around just the operation (preceded by a dollar sign) or around the entire assignment statement. Both methods produce the same results.

If you don't use the `typeset` command to declare the data type of a variable beforehand, the zsh shell attempts to automatically assign the data type. This can be dangerous when working with both integer and floating-point numbers. Take a look at this example:

```
% value=10
% value2=$(( $value1 / 3 ))
% echo $value2
3
%
```

Now that's probably not the answer you want to come out from the calculation. When you specify numbers without decimal places, the zsh shell interprets them as integer values and performs integer calculations. To ensure that the result is a floating-point number, you must specify the numbers with decimal places:

```
% value=10.
% value2=$(( $value1 / 3. ))
% echo $value2
3.33333333333333335
%
```

Now the result is in the floating-point format.

Mathematical Functions

With the zsh shell, built-in mathematical functions are either feast or famine. The default zsh shell doesn't include any special mathematical function. However, if you install the `zsh/mathfunc` module, you'll have more math functions than

you'll most likely ever need:

```
% value1=$(( sqrt(9) ))  
zsh: unknown function: sqrt  
% zmodload zsh/mathfunc  
% value1=$(( sqrt(9) ))  
% echo $value1  
3.  
%
```

That was simple! Now you have an entire math library of functions at your fingertips.

Note

There are lots of mathematical functions supported in zsh. For a complete listing of all the math functions that the `zsh/mathfunc` module provides, look at the manual page for `zshmodules`.

Structured Commands

The zsh shell provides the usual set of structured commands for your shell scripts:

- `if-then-else` statements
- `for` loops (including the C-style)
- `while` loops
- `until` loops
- `select` statements
- `case` statements

The zsh shell uses the same syntax for each of these structured commands that you're used to from the bash shell. The zsh shell also includes a different structured

command called `repeat`. The `repeat` command uses the following format:

```
repeat param
do
    commands
done
```

The `param` parameter must be a number or a mathematical operation that evaluates to a number. The `repeat` command then performs the specified commands that number of times:

```
% cat test1
#!/bin/zsh
# using the repeat command

value=$(( 10 / 2 ))
repeat $value
do
    echo "This is a test"
done
$ ./test1
This is a test
This is a test
This is a test
This is a test
This is a test
$
```

This command allows you to repeat sections of code for a set number of times based on a calculation.

Functions

The `zsh` shell supports the creation of your own functions using either the `function` command, or by defining the function name with parentheses:

```
% function functest1 {
> echo "This is the test1 function"
}
```

```

    functest2() {
> echo "This is the test2 function"
}
% functest1
This is the test1 function
% functest2
This is the test2 function
%

```

As with bash shell functions (see Chapter 16), you can define functions within your shell script and then either use global variables or pass parameters to your functions. Here's an example using a global variable:

```

% cat test3
#!/bin/zsh
# testing functions in zsh

dbl() {
    value=$(( $value * 2 ))
    return $value
}

value=10
dbl
echo The answer is $?
% ./test3
The answer is 20
%

```

You don't have to place your functions within your shell scripts. The zsh shell allows you to define your functions in separate files that it can access when trying to resolve a function name.

The zsh shell finds functions via the `fpath` environment variable. You can store your function files in any directory in this path. Here's the `fpath` value on a typical Linux workstation:

```

% echo $fpath
/usr/local/share/zsh/site-functions
/usr/share/zsh/4.2.5/functions/Completion
/usr/share/zsh/4.2.5/functions/Completion/AIX

```

```
/usr/share/zsh/4.2.5/functions/Completion/BSD
/usr/share/zsh/4.2.5/functions/Completion/Cygwin
/usr/share/zsh/4.2.5/functions/Completion/Darwin
/usr/share/zsh/4.2.5/functions/Completion/Debian
/usr/share/zsh/4.2.5/functions/Completion/Linux
/usr/share/zsh/4.2.5/functions/Completion/Mandrake
/usr/share/zsh/4.2.5/functions/Completion/Redhat
/usr/share/zsh/4.2.5/functions/Completion/Unix
/usr/share/zsh/4.2.5/functions/Completion/X
/usr/share/zsh/4.2.5/functions/Completion/Zsh
/usr/share/zsh/4.2.5/functions/MIME
/usr/share/zsh/4.2.5/functions/Misc
/usr/share/zsh/4.2.5/functions/Prompts
/usr/share/zsh/4.2.5/functions/TCP
/usr/share/zsh/4.2.5/functions/Zftp
/usr/share/zsh/4.2.5/functions/Zle
%
```

As you can see, there are lots of places the zsh shell goes hunting to resolve function names. On this system, you could place my functions in the `/usr/local/share/zsh/site-functions` directory, and the zsh shell will be able to resolve them.

However, before the zsh shell can resolve the function, you must use the `autoload` command. This command loads the function into memory for the shell to access.

Here's an example of a standalone function:

```
% cat dbl
#!/bin/zsh
# a function to double a value

dbl() {
    value=$(( $1 * 2 ))
    return $value
}

% cp dbl /usr/local/share/zsh/site-functions
%
```

Depending on how your Linux system is set up, you may have to be the root user (or use the `sudo` command) to copy files into the zsh library directories. Okay, now the function

is created in a file and stored in a directory in the `fpath`. If you try to use it, however, you'll get an error message until you load it into memory:

```
% db1 5
zsh: command not found: db1
% autoload db1
% db1 5
% echo $?
10
%
```

This also applies to shell scripts. If you have a function you need to use, you'll need to use the `autoload` command to make sure that it's available:

```
% cat test4
#!/bin/zsh
# testing an external function

autoload db1

db1 $1
echo The answer is $?
% ./test4 5
The answer is 10
%
```

Another interesting feature of the `zsh` shell is the `zcompile` command. This command processes a function file and creates a “compiled” version for the shell. This isn't really the same type of compiling you're used to in other programming languages. It does, however, put the function into a binary format that the `zsh` shell can load more quickly.

When you run the `zcompile` command, it creates a `.zwc` version of the function file. When the `autoload` command looks for the command in the `fpath`, it will see the `.zwc` version and load it instead of the text function file.

Summary

This chapter discussed two popular alternative Linux shells that you may run into. The dash shell was developed as part of the Debian Linux distribution and is mainly found in the Ubuntu Linux distribution. It's a smaller version of the Bourne shell, so it doesn't support as many features as the bash shell, which can cause problems for script writing.

The zsh shell is often found in programming environments, as it provides lots of cool features for shell script programmers. It uses loadable modules to load separate code libraries, which make using advanced functions as easy as running command line commands! There are loadable modules for lots of different functions, from complex mathematical algorithms to network applications such as FTP and HTTP.

The next section of this book dives into some specific scripting applications you might run into in the Linux environment. The next chapter shows you how to incorporate the two most popular database packages in the Linux world, MySQL and PostgreSQL, into your shell scripts for handling data.

Part IV

Advanced Shell Scripting Topics

In This Part

Chapter 23: Using a Database

Chapter 24: Using the Web

Chapter 25: Using e-mail

Chapter 26: Writing Script Utilities

Chapter 27: Advanced Shell Scripts

Chapter 23

Using a Database

In This Chapter

- Introducing MySQL
- Introducing PostgreSQL
- Creating database objects
- Writing database shell scripts

One of the problems with shell scripts is persistent data. You can store all the information you want in your shell script variables, but at the end of the script, the variables just go away. There are times when you'd like for your scripts to be able to store data that you can use later. In the old days, to store and retrieve data from a shell script required creating a file, reading data from the file, parsing the data, and then saving the data back into the file. Search for data in the file meant reading every record in the file to look for your data. Nowadays with databases being all the rage, it's a snap to interface your shell scripts with professional-quality open source databases. The two most popular open source databases used in the Linux world are MySQL and PostgreSQL. This chapter shows you how to

get these databases running on your Linux system and then spends some time getting you used to working with them from the command line. It then goes on to show you how to interact with each one using normal bash shell scripts.

The MySQL Database

By far the most popular database available in the Linux environment is the MySQL database. Its popularity has grown as a part of the Linux-Apache-MySQL-PHP (LAMP) server environment, which many Internet Web servers use for hosting online stores, blogs, and applications.

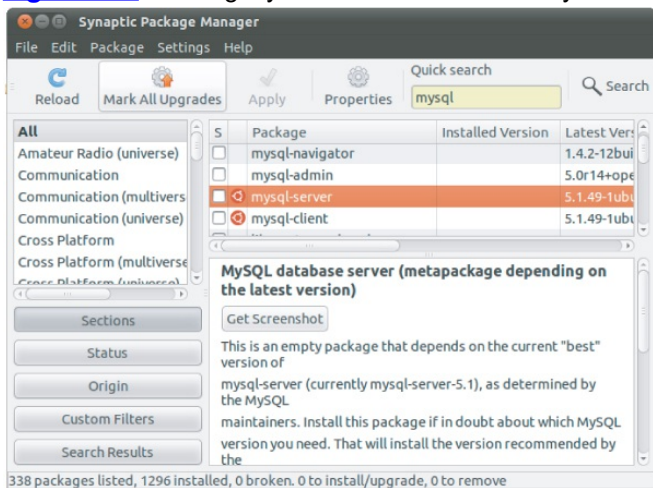
This section describes how to install a MySQL database in your Linux environment and how to create the necessary database objects to use in your shell scripts.

Installing MySQL

Installing MySQL on a Linux system used to be somewhat involved, but these days it's not uncommon for Linux distributions to use automated software installation programs (see Chapter 8). These programs not only allow you to easily download and install new software from network repositories, but also automatically check for updates for your installed software packages and install those, too.

[Figure 23.1](#) demonstrates the Add Software feature in the Ubuntu Linux distribution.

Figure 23.1 Installing MySQL on an Ubuntu Linux system

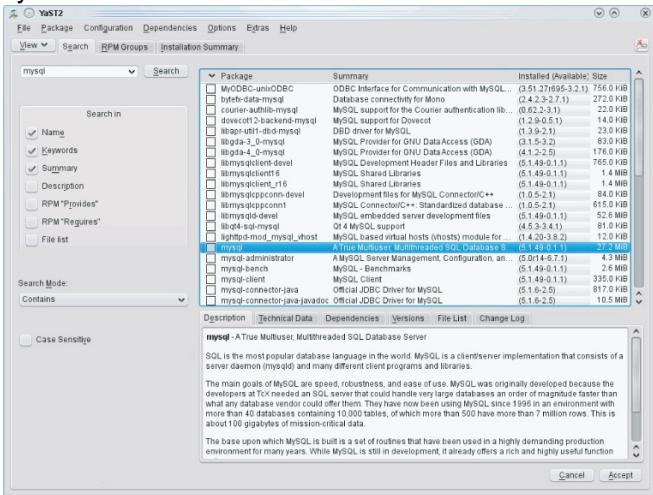


You just have to select the option for the mysql-server, and the Synaptic Package Manager downloads and installs the complete MySQL server (and client) software. It doesn't get any easier than that!

The openSUSE Linux distribution also uses an advanced software management system. [Figure 23.2](#) shows the Software Management window, where you can select packages based on software category or search for specific packages.

Figure 23.2 Installing MySQL on an openSUSE Linux

system



Again, you can select the option for the MySQL software to download and install the software packages required.

If you're using a Linux distribution that doesn't support automatic software loads, or you just want to have the latest-and-greatest version of the MySQL server, you can download an installation file directly from the MySQL Website (www.mysql.com) and install it yourself.

The MySQL Website contains both pre-built binary packages for specific Linux distributions (using either the Red Hat RPM package format or the Debian DEB

package format) and source code packages. The pre-built binary packages are easier to install, so if one is available for your Linux system, try to use it. Refer back to Chapter 8 if you need help installing either the binary or source code package for your Linux distribution.

Note

If you select to manually download and install MySQL (either the pre-built binary package or the source code download), you'll have to do some extra work configuring the MySQL server, whereas the Linux distribution installation package will install a pre-configured setup, requiring little to no work on your part!

The MySQL Client Interface

The portal to the MySQL database is the `mysql` command line interface program. This section describes the `mysql` client program and shows how to use it to interact with your database.

Connecting to the Server

The `mysql` client program allows you to connect to any MySQL database server anywhere on the network, using any user account and password. By default, if you enter the `mysql` program on a command line without any parameters, it will attempt to connect to a MySQL server running on the

same Linux system, using the Linux login username.

Most likely, this isn't how you want to connect to the database. There are lots of command line parameters that allow you to control not only which MySQL server you connect to but also the behavior of the `mysql` interface. [Table 23.1](#) shows the command line parameters you can use with the `mysql` program.

[Table 23.1](#) The `mysql` Command Line Parameters

Parameter	Description
-A	Disable automatic rehashing.
-b	Disable beep after error.
-B	Don't use a history file.
-C	Compress all information sent between the client and the server.
-D	Specify the database to use.
-e	Execute the specified statement and exit.
-E	Display query output vertically, one data field per line.
-f	Continue if an SQL error occurs.
-G	Enable named <code>mysql</code> commands.
-h	Specify the MySQL server hostname (the default is localhost).
-H	Display query output in HTML code.
-i	Ignore spaces after function names.
-N	Don't display column names in results.
-o	Ignore statements except those for the default database named on the command line.
-P	Prompt for the password for the user account.

-P	Specify the TCP port number to use for the network connection.
-q	Don't cache each query result.
-r	Display column values without escape conversion.
-s	Use silent mode.
-S	Specify a socket for connection to the localhost.
-t	Display output in table form.
-T	Display debugging information, memory, and CPU statistics when the program exits.
-u	Specify the user account to log in as.
-U	Allow only <code>UPDATE</code> and <code>DELETE</code> statements that specify key values.
-v	Use verbose mode.
-w	If the connection can't be established, wait and retry.
-X	Display query output in XHTML code.

As you can see, quite a lot of command line options are available for modifying how you log into the MySQL server.

By default, the `mysql` client program attempts to log in to the MySQL server using your Linux login name. If this name isn't configured in MySQL as a user account, you'll need to use the `-u` parameter to specify the name to log in as:

```
$ mysql -u root -p
Enter password:
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 39
Server version: 5.1.49-lubuntu8.1 (Ubuntu)

Copyright (c) 2000, 2010, Oracle and/or its affiliates. All rights reserved.
This software comes with ABSOLUTELY NO WARRANTY. This is free software,
and you are welcome to modify and redistribute it under the GPL v2 license

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.
```

mysql>

The `-p` parameter tells the `mysql` program to prompt for a password to use with the user account to log in. Once you're logged in to the server, you can start entering commands.

The mysql Commands

The `mysql` program uses two different types of commands:

- Special `mysql` commands
- Standard SQL statements

The `mysql` program uses its own set of commands that let you easily control the environment and retrieve information about the MySQL server. [Table 23.2](#) shows these commands.

Table 23.2 The mysql Commands

Command	Shortcut	Description
<code>?</code>	<code>\?</code>	Help.
<code>clear</code>	<code>\c</code>	Clear command.
<code>connect</code>	<code>\r</code>	Connect to database and server.
<code>delimiter</code>	<code>\d</code>	Set SQL statement delimiter.
<code>edit</code>	<code>\e</code>	Edit the command with the command line editor.
<code>ego</code>	<code>\G</code>	Send the command to the MySQL server and display results vertically.
<code>exit</code>	<code>\q</code>	Exit from the mysql program.
<code>go</code>	<code>\g</code>	Send command to the MySQL server.
<code>help</code>	<code>\h</code>	Display help.

nopager	\n	Disable output pager and send output to STDOUT.
note	\t	Don't send output to output file.
pager	\P	Set pager command to specified program (use more as default).
print	\p	Print current command.
prompt	\R	Change the <code>mysql</code> command prompt.
quit	\q	Quit from the <code>mysql</code> program (same as <code>exit</code>).
rehash	\#	Rebuild the command completion hash table.
source	\.	Execute the SQL script in the specified file.
status	\s	Retrieve status information from the MySQL server.
system	!	Execute a shell command on the system.
tee	\T	Append all output to the specified file.
use	\u	Use another database.
charset	\C	Change to another character set.
warnings	\W	Show warnings after every statement.
nowarning	\w	Don't show warnings after every statement.

You can use either the full command or the shortcut command directly from the `mysql` command prompt:

```
mysql> \s
```

```
-----
```

```
mysql Ver 14.12 Distrib 5.0.45, for redhat-linux-gnu (i386) using
readline 5.0
```

```
Connection id:          10
Current database:
Current user:           root@localhost
SSL:                    Not in use
Current pager:          stdout
Using outfile:          ''
Using delimiter:        ;
Server version:         5.1.49-lubuntu8.1 (Ubuntu)
Protocol version:       10
Connection:             Localhost via UNIX socket
Server characterset:    latin1
```

```

Db      characterset:  latin1
Client characterset:  latin1
Conn.   characterset:  latin1
UNIX socket:          /var/lib/mysql/mysql.sock
Uptime:               4 hours 15 min 24 sec

```

```

Threads: 1  Questions: 53  Slow queries: 0  Opens: 23  Flush tables:
1  Open tables: 17  Queries per second avg: 0.003
-----
mysql>

```

The `mysql` program implements all of the standard Structured Query Language (SQL) commands supported by the MySQL server. The “Creating MySQL Database Objects” section in this chapter discusses this in more detail.

One uncommon SQL command that the `mysql` program implements is the `SHOW` command. Using this command, you can extract information about the MySQL server, such as the databases and tables created:

```

mysql> SHOW DATABASES;
+-----+
| Database          |
+-----+
| information_schema |
| mysql              |
+-----+
2 rows in set (0.04 sec)

mysql> USE mysql;
Database changed
mysql> SHOW TABLES;
+-----+
| Tables_in_mysql   |
+-----+
| columns_priv      |
| db                 |
| func               |
| help_category     |
| help_keyword      |
| help_relation     |
| help_topic        |
| host               |

```

```
| proc |
| procs_priv |
| tables_priv |
| time_zone |
| time_zone_leap_second |
| time_zone_name |
| time_zone_transition |
| time_zone_transition_type |
| user |
+-----+
17 rows in set (0.00 sec)

mysql>
```

In this example, we used the `SHOW` SQL command to display the databases currently configured on the MySQL server, and then the `USE` SQL command to connect to a single database. Your `mysql` session can be connected to only one database at a time.

You'll notice that after each command we added a semicolon. The semicolon indicates the end of a command to the `mysql` program. If you don't use a semicolon, it prompts for more data:

```
mysql> SHOW
-> DATABASES;
+-----+
| Database |
+-----+
| information_schema |
| mysql |
+-----+
2 rows in set (0.00 sec)

mysql>
```

This feature can come in handy when you're working with long commands. You can enter part of the command on a line and press the Enter key, and then continue on the next line. This can continue for as many lines as you like until you use the semicolon to indicate the end of the command.

Note

Throughout this chapter, we use uppercase letters for SQL commands. This has become a common way to write SQL commands, but the `mysql` program allows you to specify SQL commands using either uppercase or lowercase.

Creating MySQL Database Objects

Before you can start writing your shell scripts to interact with a database, you'll need a few database objects to work with. At a minimum, you'll want to have the following:

- A unique database to store your application data
- A unique user account to access the database from your scripts
- One or more data tables to organize your data

You build all of these objects using the `mysql` program. The `mysql` program interfaces directly with the MySQL server, using SQL commands to create and modify each of the objects.

You can send any type of SQL commands to the MySQL server using the `mysql` program. This section walks you through the different SQL statements you'll need to build the basic database objects for your shell scripts.

Creating a Database

The MySQL server organizes data into *databases*. A

database usually holds the data for a single application, separating it from other applications that use the database server. Creating a separate database for each shell script application helps eliminate confusion and data mix-ups.

The SQL statement required to create a new database is:

```
CREATE DATABASE name;
```

That's pretty simple. Of course, you must have the proper privileges to create new databases on the MySQL server. The easiest way to do that is to log in as the root user account:

```
$ mysql -u root -p
Enter password:
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 39
Server version: 5.1.49-lubuntu8.1 (Ubuntu)

Copyright (c) 2000, 2010, Oracle and/or its affiliates. All rights reserved.
This software comes with ABSOLUTELY NO WARRANTY. This is free software,
and you are welcome to modify and redistribute it under the GPL v2 license

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql> CREATE DATABASE test;
Query OK, 1 row affected (0.02 sec)

mysql>
```

You can see if the new database was created by using the `SHOW` command:

```
mysql> SHOW DATABASES;
+-----+
| Database |
+-----+
| information_schema |
| mysql |
| test |
+-----+
3 rows in set (0.01 sec)
```

```
mysql>
```

Yes, it was successfully created. You should now be able to connect to the new database:

```
mysql> USE test;  
Database changed;  
mysql> SHOW TABLES;  
Empty set (0.00 sec)  
mysql>
```

The `SHOW TABLES` command allows you to see if there are any tables created. The Empty set result indicates that there aren't any tables to work with yet. Before we start creating tables, however, there's one other thing we need to do.

Creating a User Account

So far, you've seen how to connect to the MySQL server using the `root` administrator account. This account has total control over all of the MySQL server objects (much like how the `root` Linux account has complete control over the Linux system).

It's extremely dangerous to use the `root` MySQL account for normal applications. If there should be a breach of security and someone figured out the password for the root user account, all sorts of bad things could happen to your system (and data).

To prevent that, it's wise to create a separate user account in MySQL that has privileges only for the database used in the application. You do this with the `GRANT` SQL statement:


```
mysql> GRANT SELECT,INSERT,DELETE,UPDATE ON test.* TO test IDENTIFIED  
by 'test';  
Query OK, 0 rows affected (0.35 sec)  
  
mysql>
```

That's quite a long command. Let's walk through the pieces and see what it's doing.

The first section defines the privileges the user account has on what database(s). This statement allows the user account to query the database data (the select privilege), insert new data records, delete existing data records, and update existing data records.

The `test.*` entry defines the database and tables the privileges apply to. This is specified in the following format:

```
database.table
```

As you can see from this example, you're allowed to use wildcard characters when specifying the database and tables. This format applies the specified privileges to all of the tables contained in the database named test.

Finally, you specify the user account(s) to which the privileges apply. The neat thing about the `grant` command is that if the user account doesn't exist, it creates it. The `identified by` portion allows you to set a default password for the new user account.

You can test the new user account directly from the `mysql` program:

```
$ mysql test -u test -p  
Enter password:  
Welcome to the MySQL monitor.  Commands end with ; or \g.  
Your MySQL connection id is 40  
Server version: 5.1.49-lubuntu8.1 (Ubuntu)
```

```
mysql>
```

The first parameter specifies the default database to use (test), and as you've already seen, the `-u` parameter defines the user account to log in as, along with the `-p` to prompt for the password. After entering the password assigned to the test user account, you're connected to the server.

Now that you've got a database and a user account, you're ready to create some tables for the data. But first, let's take a look at the other database server you can use.

The PostgreSQL Database

The PostgreSQL database started out as an academic project, demonstrating how to incorporate advanced database techniques into a functional database server. Over the years, PostgreSQL has evolved into one of the most advanced open source database servers available for the Linux environment.

This section walks you through getting a PostgreSQL database server installed and running, and then setting up a user account and database to work with in your shell scripts.

Installing PostgreSQL

Just like MySQL, you can install the PostgreSQL database server package either by using your system's automated

software installation system, or manually by downloading it from the PostgreSQL Website (www.postgresql.org).

To use your Linux distribution's automated software installation, follow the same procedures outlined in the "Installing MySQL" section earlier in the chapter. Similar to MySQL, the PostgreSQL download area contains both pre-built packages for several Linux distributions, as well as a source code package. Again, refer back to Chapter 8 to determine if your Linux distribution uses the package management system supported by a PostgreSQL binary package. If not, you'll need to download the source code package for PostgreSQL and manually compile it on your Linux system. This can be an interesting project.

Note

Compiling a source code package requires that you have the C software development package loaded onto your Linux system. These days this is a fairly common thing for servers, but it may not be loaded if you're using a desktop Linux distribution. Please consult your individual Linux distribution documentation regarding compiling software and what software packages are required for compiling C source code projects.

Once you've installed the PostgreSQL database server, logging into the PostgreSQL server is slightly different from the MySQL server. If you remember, the MySQL server maintains its own internal database of users that can be granted access to database objects. While PostgreSQL

also has this capability, most PostgreSQL implementations (including the default source code installation) utilize the existing Linux system user accounts to authenticate PostgreSQL users.

While this can sometimes be confusing, it does make for a nice, clean way to control user accounts in PostgreSQL. You simply need to ensure that each PostgreSQL user has a valid account on the Linux system, rather than having to worry about a whole separate set of user accounts.

Another major difference for PostgreSQL is that the administrator account in PostgreSQL is called `postgres`, not `root`. Because of this requirement, when you install PostgreSQL on your Linux system, there must be a Linux user account named `postgres`.

The next section looks at how to use the `postgres` account to access the PostgreSQL server.

The PostgreSQL Command Interface

The PostgreSQL command line client program is called `psql`. This program provides complete access to the database objects configured in the PostgreSQL server. This section describes the `psql` command and shows you how to use it to interact with your PostgreSQL server.

Connecting to the Server

The `psql` client program provides the command line interface to the PostgreSQL server. As you would expect, it

uses command line parameters to control what features are enabled in the client interface. Each option uses either a long or short name format. [Table 23.3](#) shows the command line parameters available.

Table 23.3 The psql Command Line Parameters

Short Name	Long Name	Description
-a	--echo-all	Display all SQL lines processed from a script file in the output.
-A	--no-align	Set the output format to unaligned mode. Data is not displayed as a formatted table.
-c	--command	Execute the specified SQL statement and exit.
-d	--dbname	Specify the database to connect with.
-e	--echo-queries	Echo all queries to the screen.
-E	--echo-hidden	Echo hidden psql meta-commands to the screen.
-f	--file	Execute SQL commands from the specified file and exit.
-F	--field-separator	Specify the character used to separate column data when in unaligned mode. The default is a comma.
-h	--host	Specify the IP address or hostname of the remote PostgreSQL server.
-l	--list	Display a list of available databases on the server and exit.
-o	--output	Redirect query output to the specified file.
-p	--port	Specify the PostgreSQL server TCP port to connect with.
-P	--pset	Set the table printing option specified to a specified value.
-q	--quiet	Quiet mode, doesn't display output messages.

-R	--record-separator	Use the specified character as the record separator. The default is the newline character.
-s	--single-step	Prompt to continue or cancel after every SQL query.
-S	--single-line	Specify that the Enter key defines the end of a SQL query instead of a semicolon.
-t	--tuples-only	Disables column headers and footers in table output.
-T	--table-attr	Use the HTML table tag specified when in HTML mode.
-U	--username	Use the specified username to connect to the PostgreSQL server.
-v	--variable	Set the specified variable to a specified value.
-V	--version	Display the psql version number and exit.
-W	--password	Force a password prompt.
-x	--expanded	Enable expanded table output to display additional information for records.
-X	--nopsqlrc	Don't process the psql startup file.
-?	--help	Display the psql command line help and exit.

As mentioned in the previous section, the administrative account for PostgreSQL is called `postgres`. Because PostgreSQL uses Linux user accounts to validate users, you must be logged in as the `postgres` Linux account to access the PostgreSQL server as the `postgres` user.

To get around this problem, you can use the `sudo` command to run the `psql` command line program as the `postgres` user account.

```
$ sudo -u postgres psql
[sudo password for rich]:
psql (8.4.5)
Type "help" for help.

postgres=#
```

The default `psql` prompt indicates the database you are connected to. The pound sign (#) in the prompt indicates that you're logged in with the administrative user account. You're now ready to start entering some commands to interact with the PostgreSQL server.

The psql Commands

Similarly to the `mysql` program, the `psql` program uses two different types of commands:

- Standard SQL statements
- PostgreSQL meta-commands

PostgreSQL meta-commands allow you to easily extract information about the database environment, plus set features for the `psql` session. A meta-command is indicated by using a backslash. There are lots of PostgreSQL meta-commands for lots of different settings and features, but there's no reason to start worrying about them all right away. The most commonly used ones are:

- `\l` to list the available databases
- `\c` to connect to a database
- `\dt` to list the tables in a database
- `\du` to list the PostgreSQL users
- `\z` to list table privileges

- `\?` to list all of the available meta-commands
- `\h` to list all of the available SQL commands
- `\q` to exit the database

If you ever need to find a meta-command, just enter the `\?` meta-command. You'll see a list, along with an explanation, of all the available meta-commands.

To test the meta-commands, use the `\l` meta-command to list the available databases:

```
postgres=# \l
List of databases

```

Name	Owner	Encoding	Collation	Ctype	Access
postgres	postgres	UTF8	en_US.utf8	en_US.utf8	
template0	postgres	UTF8	en_US.utf8	en_US.utf8	=c/postgres
template1	postgres	UTF8	en_US.utf8	en_US.utf8	=c/postgres

```
(3 rows)
postgres=#
```

The listing shows the available databases on the server, along with their features (we truncated the Access column entries so they would fit on the page). These are the default databases provided by the PostgreSQL server. The `postgres` database maintains all of the system data for the server. The `template0` and `template1` databases provide default database templates for you to copy when creating a new database.

You're now ready to start working on your own data in PostgreSQL.

Creating PostgreSQL Database Objects

This section walks you through the process of creating your

database and a user account to access it. You'll see that while some of the work in PostgreSQL is exactly the same as in MySQL, some of it is completely different.

Creating a Database Object

Creating a database is one of those actions that's the same as in MySQL. Remember to be logged in as the `postgres` administrative account to create the new database:

```
$ sudo -u postgres psql
psql (8.4.5)
Type "help" for help.

postgres=# CREATE DATABASE test;
CREATE DATABASE
postgres=#
```

After you create the database, use the `\l` meta-command to see if it appears in the listing, then the `\c` meta-command to connect to it:

```
postgres=# \l
List of databases

```

Name	Owner	Encoding	Collation	Ctype	Access
postgres	postgres	UTF8	en_US.utf8	en_US.utf8	
template0	postgres	UTF8	en_US.utf8	en_US.utf8	=c/postgres
template1	postgres	UTF8	en_US.utf8	en_US.utf8	=c/postgres
test	postgres	UTF8	en_US.utf8	en_US.utf8	

```
(4 rows)

postgres=# \c test
psql (8.4.5)
You are now connected to database "test".
test=#
```

When you connect to the `test` database, the `psql` prompt changes to indicate the new database name. This is a great reminder when you're ready to create your database objects, so you can easily tell where you are in the system.

Note

PostgreSQL adds another layer of control to the database called the *schema*. A database can contain multiple schemas, each schema containing multiple tables. This allows you to subdivide a database for specific applications or users.

By default, every database contains one schema, called *public*. If you're only going to have one application use the database, you're fine with just using the public schema. If you'd like to really get fancy, you can create new schemas. For this example, use the public schema for tables.

Creating User Accounts

After creating the new database, the next step is to create a user account that has access to it for your shell scripts. As you've already seen, user accounts in PostgreSQL are significantly different from those in MySQL.

User accounts in PostgreSQL are called *Login Roles*. The PostgreSQL server matches Login Roles to the Linux system user accounts. Because of this, there are two common approaches to creating Login Roles to run shell scripts that access the PostgreSQL database:

- Create a special Linux account with a matching PostgreSQL Login Role to run all your shell scripts.
- Create PostgreSQL accounts for each Linux user account that needs to run shell scripts to access the database.

For this example, let's choose the second method and create a PostgreSQL account that matches our Linux system account. This way, you can run shell scripts that access the PostgreSQL database directly from your Linux user account.

First, you must create the Login Role:

```
test=# CREATE ROLE rich login;  
CREATE ROLE  
test=#
```

That was simple enough. Without the `login` parameter, the role is not allowed to log in to the PostgreSQL server, but it can be assigned privileges. This type of role is called a *group role*. Group roles are great if you're working in a large environment with lots of users and tables. Instead of having to keep track of which user has which type of privileges for which tables, you just create group roles for specific types of access to tables, and then assign the login roles to the proper group role.

For simple shell scripting, you most likely won't need to worry about creating group roles, and just assign privileges directly to the login roles. That's what we'll do in this example.

However, PostgreSQL handles privileges a bit differently than MySQL. It doesn't allow you to grant overall privileges to all objects in a database that filter down to the table level. Instead, you'll need to grant privileges for each individual table you create. While this is somewhat of a pain, it certainly helps enforce strict security policies. You'll have to hold off assigning privileges until you've created a table.

That's the next step in our process.

Working with Tables

Now that you've got your MySQL or PostgreSQL server running and created a new database and a user account for accessing it, it's time to start working with data! Fortunately, both the `mysql` and `psql` programs use standard SQL to create and manage data tables. This section walks you through the SQL required to create tables, insert and remove data, and query existing data in both environments.

Creating a Table

Both the MySQL and PostgreSQL servers are considered *relational* databases. In a relational database, data is organized by *data fields*, *records*, and *tables*. A data field is a single piece of information, such as an employee's last name or a salary. A record is a collection of related data fields, such as the employee ID number, last name, first name, address, and salary. Each record indicates one set of the data fields.

The table contains all of the records that hold the related data. Thus, you'll have a table called `Employees` that holds the records for each employee.

To create a new table in the database, you need to use the `CREATE TABLE` SQL command:

```
$ mysql test -u root -p
```

```
Enter password:
mysql> CREATE TABLE employees (
    -> empid int not null,
    -> lastname varchar(30),
    -> firstname varchar(30),
    -> salary float,
    -> primary key (empid));
Query OK, 0 rows affected (0.14 sec)

mysql>
```

First, notice that to create the new table we needed to log in to MySQL using the `root` user account because the `test` user doesn't have privileges to create a new table. Next, notice that we specified the `test` database on the `mysql` program command line. If we hadn't done that, we would need to use the `USE` SQL command to connect to the `test` database.

Caution

It's extremely important that you make sure you're in the right database before creating the new table. Also, make sure that you're logged in using the administrative user account (`root` for MySQL and `postgres` for PostgreSQL) to create the tables.

Each data field in the table is defined using a data type. The MySQL and PostgreSQL databases support lots of different data types. [Table 23.4](#) shows some of the more popular data types you may need.

Table 23.4 MySQL and PostgreSQL Data Types

Data Type	Description
char	A fixed-length string value

varchar	A variable-length string value
int	An integer value
float	A floating-point value
Boolean	A Boolean true/false value
Date	A date value in YYYY-MM-DD format
Time	A time value in HH:mm:ss format
Timestamp	A date and time value together
Text	A long string value
BLOB	A large binary value, such as an image or video clip

The `empid` data field also specifies a *data constraint*. A data constraint restricts what type of data you can enter to create a valid record. The `not null` data constraint indicates that every record must have an `empid` value specified.

Finally, the `primary key` defines a data field that uniquely identifies each individual record. This means that each data record must have a unique `empid` value in the table.

After creating the new table, you can use the appropriate command to ensure that it's created. In `mysql`, it's the `show tables` command:

```
mysql> show tables;
+-----+
| Tables_in_test |
+-----+
| employees      |
+-----+
1 row in set (0.00 sec)

mysql>
```

And in `psql` it's the `\dt` meta-command:

```
test=# \dt
List of relations
```

Schema	Name	Type	Owner
public	employees	table	postgres

(1 row)

```
test=#
```

As you may remember from the “Creating PostgreSQL Database Objects” section earlier in this chapter, in PostgreSQL you need to assign privileges at the table level. Now that you have a table, you'll need to give your Login Role access to it:

```
$ sudo -u postgres psql
psql (8.4.5)
Type "help" for help.

postgres=# \c test
psql (8.4.5)
You are now connected to database "test".
test=# GRANT SELECT, INSERT, DELETE, UPDATE ON public.employees TO rich;
GRANT
test=#
```

The format to specify the table must include the schema name, which by default is `public`. Also, remember to perform this command as the `postgres` login role, and to connect to the test database.

With the table created, you're now ready to start saving some data. The next section covers how to do that.

Inserting and Deleting Data

Not surprisingly, you use the `INSERT` SQL command to insert new data records into the table. Each `INSERT` command must specify the data field values for the MySQL or PostgreSQL server to accept the record.

The format of the `INSERT SQL` command is:

```
INSERT INTO table VALUES (...)
```

The values are a comma-separated list of the data values for each data field:

```
$ mysql test -u test -p
Enter password:
```

```
mysql> INSERT INTO employees VALUES (1, 'Blum', 'Rich', 25000.00);
Query OK, 1 row affected (0.35 sec)
```

or, in PostgreSQL:

```
$ psql test
psql (8.4.5)
Type "help" for help.
```

```
test=> INSERT INTO employees VALUES (1, 'Blum', 'Rich', 25000.00);
INSERT 0 1
test=>
```

The MySQL example uses the `-u` command line prompt to log in as the `test` user account that was created in MySQL. The PostgreSQL example uses the current Linux user account, so it just uses the `rich` user account that was created.

The `INSERT` command pushes the data values you specify into the data fields in the table. If you attempt to add another record that duplicates the `empid` data field value, you'll get an error message:

```
mysql> INSERT INTO employees VALUES (1, 'Blum', 'Barbara', 45000.00);
ERROR 1062 (23000): Duplicate entry '1' for key 1
```

However, if you change the `empid` value to a unique value, everything should be okay:

```
mysql> INSERT INTO employees VALUES (2, 'Blum', 'Barbara', 45000.00);
Query OK, 1 row affected (0.00 sec)
```

You should now have two data records in your table.

If you need to remove data from your table, you use the `DELETE` SQL command. However, you need to be very careful with it.

The basic `DELETE` command format is:

```
DELETE FROM table;
```

where *table* specifies the table to delete records from. There's just one small problem with this command: It removes all of the records in the table.

To just specify a single record or a group of records to delete, you must use the `WHERE` clause. The `WHERE` clause allows you to create a filter that identifies which records to remove. You use the `WHERE` clause like this:

```
DELETE FROM employees WHERE empid = 2;
```

This restricts the deletion process to all of the records that have an `empid` value of 2. When you execute this command, the `mysql` program returns a message indicating how many records matched the filter:

```
mysql> DELETE FROM employees WHERE empid = 2;  
Query OK, 1 row affected (0.29 sec)
```

As expected, only one record matched the filter and was removed.

Querying Data

Once you've got all of your data in your database, it's time to start running reports to extract information.

The workhorse for all of your querying is the SQL `SELECT` command. The `SELECT` command is extremely versatile, but

with versatility comes complexity.

The basic format of a `SELECT` statement is:

```
SELECT datafields FROM table
```

The `datafields` parameter is a comma-separated list of the data field names you want the query to return. If you want to receive all of the data field values, you can use an asterisk as a wildcard character.

You must also specify the specific table you want the query to search. To get meaningful results, you must match your query data fields with the proper table.

By default, the `SELECT` command returns all of the data records in the specified table:

```
mysql> SELECT * FROM employees;
+-----+-----+-----+-----+
| empid | lastname | firstname | salary |
+-----+-----+-----+-----+
| 1 | Blum | Rich | 25000 |
| 2 | Blum | Barbara | 45000 |
| 3 | Blum | Katie Jane | 34500 |
| 4 | Blum | Jessica | 52340 |
+-----+-----+-----+-----+
4 rows in set (0.00 sec)

mysql>
```

You can use one or more modifiers to define how the database server returns the data requested by the query. Here's a list of commonly used modifiers:

- **WHERE:** Displays a subset of records that meet a specific condition
- **ORDER BY:** Displays records in a specified order
- **LIMIT:** Displays only a subset of records

The `WHERE` clause is the most common `SELECT` command

modifier. It allows you to specify conditions to filter data from the result set. Here's an example of using the `WHERE` clause:

```
mysql> SELECT * FROM employees WHERE salary > 40000;
```

```
+-----+-----+-----+-----+
| empid | lastname | firstname | salary |
+-----+-----+-----+-----+
|      2 | Blum     | Barbara   | 45000  |
|      4 | Blum     | Jessica   | 52340  |
+-----+-----+-----+-----+
```

```
2 rows in set (0.01 sec)
```

```
mysql>
```

Now you can see the power of adding database access to your shell scripts! You can easily control your data management needs just with a few SQL commands and the `mysql` or `psql` programs. The next section describes how you can incorporate these features into your shell scripts.

Using the Database in Your Scripts

Now that you've got a working database going, it's finally time to turn our attention back to the shell scripting world. This section describes what you need to do to interact with your databases using shell scripts.

Connecting to the Databases

Obviously, to connect to the databases you'll need to somehow utilize the `mysql` or `psql` programs in your shell

script. This isn't too complicated a process, but there are a few things you'll need to watch out for.

Finding the Programs

The first hurdle you'll need to complete is to figure out just where the `mysql` and `psql` command line client programs are on the Linux system. The one downside to Linux software installs is that often different Linux distributions place software packages in different locations.

Fortunately, there's the `which` command. The `which` command tells you where the shell would find a command if it attempted to run it from the command line:

```
$ which mysql
/usr/bin/mysql
$ which psql
/usr/bin/psql
$
```

The easiest way to handle this information is to assign it to an environment variable and then use that in your shell script when you want to reference the appropriate program:

```
MYSQL='which mysql'
PSQL='which psql'
```

Now the `$MYSQL` variable points to the executable for the `mysql` program, and the `$PSQL` variable points to the executable for the `psql` program.

Logging into the Server

After finding the location of the client programs, you can use them in your scripts to access the database servers.

For the PostgreSQL server, this is easy:

```
$ cat ptest1
#!/bin/bash
# test connecting to the PostgreSQL server

PSQL='which psql'

$PSQL test
$ ./ptest1
psql (8.4.5)
Type "help" for help.

test=>
```

Because the script is being run from a Linux user account, and a PostgreSQL account exists with the same name, all you need to specify on the `psql` command line is the name of the database to connect with. The `ptest1` script connected to the `test` database and left you at the `psql` prompt inside that database.

If you've created a special user account in MySQL for your shell scripts, you'll need to specify that on the `mysql` command line:

```
$ cat mtest1
#!/bin/bash
# test connecting to the MySQL server

MYSQL='which mysql'

$MYSQL test -u test -p
$ ./mtest1
Enter password:
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 43
Server version: 5.1.49-lubuntu8.1 (Ubuntu)
```

Copyright (c) 2000, 2010, Oracle and/or its affiliates. All rights reserved.
This software comes with ABSOLUTELY NO WARRANTY. This is free software,
and you are welcome to modify and redistribute it under the GPL v2 license

Type `'help;'` or `'\h'` for help. Type `'\c'` to clear the current input statement.

```
mysql>
```

That worked but not all that well for a non-interactive script. The `-p` command line parameter caused `mysql` to pause and ask for the password. You can solve that problem by including the password on the command line:

```
$MYSQL test -u test -ptest
```

This, however, is not a good idea. Anyone who has access to your script will know the user account and password for your database.

To solve this problem, you can use a special configuration file used by the `mysql` program. The `mysql` program uses the `$HOME/.my.cnf` file to read special startup commands and settings. One of those settings is the default password for `mysql` sessions started by the user account.

To set the default password in this file, just create the following:

```
$ cat .my.cnf
[client]
password = test
$ chmod 400 .my.cnf
$
```

The `chmod` command is used to restrict the `.my.cnf` file so that only you can view it. You can test this now from the command line:

```
$ mysql test -u test
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Welcome to the MySQL monitor.  Commands end with ; or \g.
```

Your MySQL connection id is 44
Server version: 5.1.49-lubuntu8.1 (Ubuntu)

Copyright (c) 2000, 2010, Oracle and/or its affiliates. All rights reserved.
This software comes with ABSOLUTELY NO WARRANTY. This is free software,
and you are welcome to modify and redistribute it under the GPL v2 license

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql>

Perfect! Now you don't have to include the password on the command line in your shell scripts.

Sending Commands to the Server

After establishing the connection to the server, you'll want to send commands to interact with your database. There are two methods to do this:

- Send a single command and exit.
- Send multiple commands.

To send a single command, you must include the command as part of the `mysql` or `psql` command line.

For the `mysql` command, you do this using the `-e` parameter:

```
$ cat mtest2
#!/bin/bash
# send a command to the MySQL server

MYSQL='which mysql'

$MYSQL test -u test -e 'select * from employees'
$ ./mtest2
```

empid	lastname	firstname	salary
1	Blum	Rich	25000
2	Blum	Barbara	45000
3	Blum	Katie Jane	34500

```

4 | Blum          | Jessica      | 52340 |
+-----+-----+-----+
$

```

For the `psql` command, you do this using the `-c` parameter:

```

$ cat ptest2
#!/bin/bash
# send a command to the PostgreSQL server

PSQL='which psql'

$PSQL test -c 'select * from employees'
$ ./ptest2

```

empid	lastname	firstname	salary
1	Blum	Rich	25000
2	Blum	Barbara	45000
3	Blum	Katie Jane	34500
4	Blum	Jessica	52340

```

(4 rows)
$

```

The database servers return the results from the SQL commands to the shell scripts, which display them in `STDOUT`.

If you need to send more than one SQL command, you can use file redirection (see Chapter 14). To redirect lines in the shell script, you must define an *end of file* string. The end of file string indicates the beginning and end of the redirected data.

The following is an example of defining an end of file string, with data in it:

```

$ cat mtest3
#!/bin/bash
# sending multiple commands to MySQL

MYSQL='which mysql'
$MYSQL test -u test <<EOF
show tables;
select * from employees where salary > 40000;
EOF
$ ./mtest3

```



```
Tables_in_test
employees
empid      lastname      firstname      salary
2          Blum          Barbara        45000
4          Blum          Jessica        52340
$
```

The shell redirects everything with the `EOF` delimiters to the `mysql` command, which executes the lines as if you typed them yourself at the prompt. Using this method, you can send as many commands to the MySQL server as you need. You'll notice, however, that there's no separation between the output from each command. In the next section, "Formatting Data," you'll see how to fix this problem.

Note

You should also notice that the `mysql` program changed the default output style when you used the redirected input method. Instead of creating the ASCII symbol boxes around the data, the `mysql` program detected that the input was redirected, so it returned just the raw data. This will come in handy when you need to extract the individual data elements.

The same technique also works for the `psql` program:

```
$ cat ptest3
#!/bin/bash
# sending multiple commands to PostgreSQL

PSQL='which psql'

$PSQL test <<EOF
\dt
select * from employees where salary > 40000;
```

```

EOF
$ ./ptest3
          List of relations
 Schema |      Name      | Type | Owner
-----+-----+-----+-----
 public | employees | table | postgres
(1 row)

 empid | lastname | firstname | salary
-----+-----+-----+-----
      2 | Blum      | Barbara   | 45000
      4 | Blum      | Jessica   | 52340
(2 rows)
$

```

The `psql` program displays the output from each command directly to `STDOUT` in the order in which you specify them.

Of course, you're not limited to just retrieving data from the tables. You can use any type of SQL command in your script, such as an `INSERT` statement:

```

$ cat mtest4
#!/bin/bash
# send data to the table in the MySQL database

MYSQL='which mysql'

if [ $# -ne 4 ]
then
    echo "Usage: mtest4 empid lastname firstname salary"
else
    statement="INSERT INTO employees VALUES ($1, '$2', '$3', $4)"
    $MYSQL test -u test << EOF
    $statement
EOF
if [ $? -eq 0 ]
then
    echo Data successfully added
else
    echo Problem adding data
fi
fi

$ ./mtest4
Usage: mtest4 empid lastname firstname salary
$ ./mtest4 5 Blum Jasper 100000

```

```
Data added successfully
$
$ ./mtest4 5 Blum Jasper 100000
ERROR 1062 (23000) at line 1: Duplicate entry '5' for key 1
Problem adding data
$
```

This example demonstrates a few things about using this technique. When you specify the end of file string, it must be the only thing on the line, and the line must start with the string. If we had indented the `EOF` text to match the rest of the if-then indentation, it wouldn't work.

Inside the `INSERT` statement, notice that there are single quotes around the text values and double quotes around the entire `INSERT` statement. It's important not to mix up the quotes used for the string values with the quotes used to define the script variable text.

Also, notice how we used the special `$_` variable to test the exit status of the `mysql` program. This helps you determine whether the command failed or not.

Just sending output from the commands to `STDOUT` is not the easiest way to manage and manipulate the data. The next section shows you some tricks you can use to help your scripts capture data retrieved from the database.

Formatting Data

The standard output from the `mysql` and `psql` commands doesn't lend itself too well for data retrieval. If you need to actually do something with the data you retrieve, you'll need to do some fancy data manipulation. This section describes some of the tricks you can use to help extract

data from your database reports.

Assigning Output to a Variable

The first step in trying to capture database data is to redirect the output from the `mysql` and `psql` commands in an environment variable. This allows you to use the output information in other commands. Here's an example:

```
$ cat mtest5
#!/bin/bash
# redirecting SQL output to a variable

MYSQL='which mysql'

dbs='$MYSQL test -u test -Bse `show databases`'
for db in $dbs
do
    echo $db
done
$ ./mtest5
information_schema
test
$
```

This example uses two additional parameters on the `mysql` program command line. The `-B` parameter specifies for the `mysql` program to work in batch mode, and in combination with the `-s` (silent) parameter, the column headings and formatting symbols are suppressed.

By redirecting the output of the `mysql` command to a variable, this example is able to step through the individual values of each returned record.

Using Formatting Tags

In the previous example, you saw how adding the `-B` and `-s`

parameters to the `mysql` program command line allows you to suppress the output heading information, so all you get is data. There are a few other parameters that you can use to help make life easier for you.

Generating data for Web pages is a popular thing these days. Both the `mysql` and `psql` programs provide an option to display the output using HTML format. In both, this is enabled using the `-H` command line parameter:

```
$ psql test -H -c 'select * from employees where empid = 1'
<table border="1">
<tr>
  <th align="center">empid</th>
  <th align="center">lastname</th>
  <th align="center">firstname</th>
  <th align="center">salary</th>
</tr>
<tr valign="top">
  <td align="right">1</td>
  <td align="left">Blum</td>
  <td align="left">Rich</td>
  <td align="right">25000</td>
</tr>
</table>
<p>(1 row)<br />
</p>
$
```

The `mysql` program also supports an additional popular format, called the Extensible Markup Language (XML). This language uses HTML-like tags to identify data names and values.

For the `mysql` program, you do this using the `-X` command line parameter:

```
$ mysql test -u test -X -e 'select * from employees where empid = 1'
<?xml version="1.0"?>

<resultset statement="select * from employees">
<row>
  <field name="empid">1</field>
```

```
<field name="lastname">Blum</field>
<field name="firstname">Rich</field>
<field name="salary">25000</field>
</row>
</resultset>
$
```

Using XML, you can easily identify individual rows of data, along with the individual data values in each record.

Summary

This chapter discussed the ability to save, modify, and retrieve data from your shell scripts in databases. You can easily access both the MySQL and PostgreSQL database servers directly from your shell scripts.

After installing the MySQL and PostgreSQL servers, you can use their respective client programs to access the servers from the command line, or from shell scripts. The `mysql` client program provides the command line interface to the MySQL server. You can send SQL commands as well as customized MySQL commands to the server from your shell scripts and then retrieve the results.

The `psql` client program operates the same way for the PostgreSQL server. There are plenty of command line parameters you can use to help format your data in just the right way.

Both client programs allow you to send either a single command to the server or use input redirection to send a batch of commands. The programs normally send the output data from the server to `STDOUT`, but you can redirect

the output to a variable and use that information in your shell script.

The next chapter examines the World Wide Web. Getting your shell scripts to interface with Websites on the Internet is a tricky task, but once you master it the world is yours from which to retrieve data.

Chapter 24

Using the Web

In This Chapter

- Surfing with Lynx
- Exploring with cURL
- Client/server programming in zsh

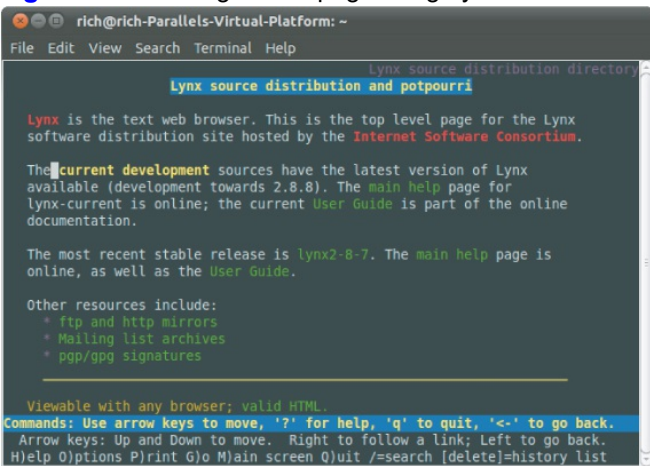
Often when you think of shell script programming, the last thing you think of is the Internet. The command line world often seems foreign to the fancy, graphical world of the Internet. There are, however, several different utilities you can easily use in your shell scripts to gain access to data content on the Web, as well as on other network devices. This chapter walks you through three popular methods for getting your shell scripts to interact with the network world.

The Lynx Program

Almost as old as the Internet itself, the Lynx program was created in 1992 by students at the University of Kansas as a text-based browser. Because it's text-based, the Lynx

program allows you to browse Websites directly from a terminal session, replacing the fancy graphics on Web pages with HTML text tags. This allows you to surf the Internet from just about any type of Linux terminal. A sample Lynx screen is shown in [Figure 24.1](#).

Figure 24.1 Viewing a Web page using Lynx



```
rich@rich-Parallels-Virtual-Platform: ~  
File Edit View Search Terminal Help  
Lynx source distribution directory  
Lynx source distribution and potpourri  
Lynx is the text web browser. This is the top level page for the Lynx  
software distribution site hosted by the Internet Software Consortium.  
The current development sources have the latest version of Lynx  
available (development towards 2.8.8). The main help page for  
lynx-current is online; the current User Guide is part of the online  
documentation.  
The most recent stable release is lynx2-8-7. The main help page is  
online, as well as the User Guide.  
Other resources include:  
* ftp and http mirrors  
* Mailing list archives  
* pgp/gpg signatures  
Viewable with any browser; valid HTML.  
Commands: Use arrow keys to move, '?' for help, 'q' to quit, '<-' to go back.  
Arrow keys: Up and Down to move. Right to follow a link; Left to go back.  
H)elp O)ptions P)rint G)o M)ain screen Q)uit /=search [delete]=history list
```

Lynx uses the standard keyboard keys to navigate around the Web page. Links appear as highlighted text within the Web page. Using the right arrow key allows you to follow a link to the next Web page.

You may be wondering how you can use a graphical text program in your shell scripts. The Lynx program also

provides a feature that allows you to dump the text contents of a Web page to `STDOUT`. This feature is great for mining for data contained within a Web page. This section describes how to use the Lynx program within your shell scripts to extract data from Internet Websites.

Installing Lynx

Even though the Lynx program is somewhat old, it's still in active development. At the time of this writing, the latest version of Lynx is version 2.8.8, released in June 2010, with a new release in development. Because of its popularity among shell script programmers, many Linux distributions install the Lynx program in their default installations.

If you're using an installation that doesn't provide the Lynx program, check your distribution's installation packages. Most likely you'll find it there for easy installation.

If your distribution doesn't include the Lynx package, or if you just want the latest version, you can download the source code from the lynx.isc.org website and compile it yourself (assuming that you've got the C development libraries installed on your Linux system). See Chapter 8 for information on how to compile and install source code distribution packages.

Note

The Lynx program uses the curses text-graphics library in Linux. Most distributions have this installed

by default. If your distribution doesn't, consult your particular distribution's instructions on installing the curses library before trying to compile Lynx.

The next section describes how to use the `lynx` command from the command line.

The lynx Command Line

The `lynx` command line command is extremely versatile in what information it can retrieve from the remote Website. When you view a Web page in your browser, you're only seeing part of the information that's transferred to your browser. Web pages consist of three types of data elements:

- HTTP headers
- Cookies
- HTML content

HTTP headers provide information about the type of data sent in the connection, the server sending the data, and the type of security used in the connection. If you're sending special types of data, such as video or audio clips, the server identifies that in the HTTP headers. The Lynx program allows you to view all of the HTTP headers sent within a Web page session.

If you've done any type of Web browsing, no doubt you're familiar with Web page *cookies*. Websites use cookies to store data about your Website visit for future use. Each individual site can store information, but only access the

information that it sets. The `lynx` command provides options for you to view cookies sent by Web servers, as well as reject or accept specific cookies sent from servers.

The Lynx program allows you to view the actual HTML content of the Web page in three different formats:

- In a text-graphics display on the terminal session using the curses graphical library
- As a text file, dumping the raw data from the Web page
- As a text file, dumping the raw HTML source code from the Web page

For shell scripts, viewing the raw data or HTML source code is a gold mine. Once you capture the data retrieved from a Website, you can easily extract individual pieces of information.

As you can see, the Lynx program is extremely versatile in what it can do. However, with versatility comes complexity, especially when it comes to command line parameters. The Lynx program is one of the more complex programs you'll run into in the Linux world.

The basic format of the `lynx` command is:

```
lynx options URL
```

where `URL` is the HTTP or HTTPS destination you want to connect to, and `options` are one or more options that modify the behavior of Lynx as it interacts with the remote Website. There are options for just about any type of Web interaction required by Lynx. [Table 24.1](#) shows all of the available command line parameters you can use with the `lynx`

command.

Table 24.1 The lynx Command Parameters

Parameter	Description
-	Receive options and arguments from STDIN.
-accept_all_cookies	Accept cookies without prompting if <code>Set-Cookie</code> handling is on. Set to <code>off</code> by default.
-anonymous	Apply restrictions for anonymous account.
-assume_charset=name	Default charset for documents that don't specify one.
-assume_local_charset=name	Default charset for local files.
-assume_unrec_charset=name	Default charset to use instead of unrecognized charsets.
-auth=id:pw	Authentication information for protected documents.
-base	Prepend a request URL comment and BASE tag to text/HTML outputs for <code>-source</code> dumps.
-bibhost=URL	Local bibp server URL (default http://bibhost/).
-book	Use the bookmark page as the start page. Set to <code>off</code> by default.
-buried_news	Toggle scanning of news articles for buried references. Set to <code>on</code> by default.
-cache=n	Number of documents cached in memory.
-case	Enable case-sensitive user searching. Set to <code>off</code> by default.
-center	Toggle center alignment in HTML <code><table></code> tags. Set to <code>off</code> by default.
-cfg=filename	Specify a configuration file other than the default <code>lynx.cfg</code> file.

-child	Exit on left arrow in startfile, and disable save to disk.
-cmd_log=filename	Log keystroke commands to the specified file.
-cmd_script=filename	Read keystroke commands from the specified file.
-connect_timeout=n	Set the connection timeout (in seconds). The default is 18,000 seconds.
-cookie_file=filename	Specify the file to use to read cookies.
-cookie_save_file=filename	Specify the file to use to store cookies.
-cookies	Toggle handling of <code>Set-Cookie</code> headers. Set to <code>on</code> by default.
-core	Toggle forced core dumps on fatal errors. Set to <code>off</code> by default.
-crawl	With <code>-traversal</code> , output each page to a file with <code>-dump</code> , format output as with <code>-traversal</code> , but to <code>STDOUT</code> .
-curses_pads	Use curses pad feature to support left/right shifting. Set to <code>on</code> by default.
-debug_partial	Display incremental display stages with <code>MessageSecs</code> delay. Set to <code>off</code> by default.
-delay=n	Set delay at <code>statusline</code> message (in seconds). Set to <code>0.000</code> by default.
-display=display	Set the display variable for XWindow programs.
-display_charset=name	The charset for the terminal output.
-dont_wrap_pre	Don't wrap text in <code><pre></code> sections when <code>-dump</code> and <code>-crawl</code> set. Mark wrapped lines in interactive session. Set to <code>on</code> by default.
-dump	Dump the first URL to <code>STDOUT</code> and exit.
-editor=editor	Enable edit mode with the specified editor.

<code>-emacskeys</code>	Enable emacs-like key movement. Set to <code>off</code> by default.
<code>-enable_scrollback</code>	Toggle compatibility with scrollbar keys. Set to <code>off</code> by default.
<code>-error_file=filename</code>	Write the HTTP status code to the specified file.
<code>-exec</code>	Enable local program execution.
<code>-force_empty_hrefless_a</code>	Force <code><a></code> elements without an href attribute to be empty. Set to <code>off</code> by default.
<code>-force_html</code>	Force the first document to be interpreted as HTML. Set to <code>off</code> by default.
<code>-force_secure</code>	Require the secure flag for SSL cookies. Set to <code>off</code> by default.
<code>-forms_options</code>	Use forms-based options menu. Set to <code>on</code> by default.
<code>-from</code>	Enable transmission of From headers. Set to <code>on</code> by default.
<code>-ftp</code>	Disable FTP access. Set to <code>off</code> by default.
<code>-get_data</code>	Read data for get forms from STDIN, terminated by <code>---</code> .
<code>-head</code>	Send a <code>HEAD</code> request. Set to <code>off</code> by default.
<code>-help</code>	Print usage message.
<code>-hiddenlinks=option</code>	Specify how to handle hidden links. <i>Option</i> can be <code>merge</code> , <code>listonly</code> , or <code>ignore</code> .
<code>-historical</code>	Use <code>></code> instead of <code>--></code> as terminator for comments. Set to <code>off</code> by default.
<code>-homepage=URL</code>	Set home page separate from start page.
<code>-image_links</code>	Enable inclusion of links for all images. Set to <code>off</code> by default.
<code>-index=URL</code>	Set the default index file name.
	Include ISMAP links when client-side MAPs are

-ismap	present. Set to <code>off</code> by default.
-link= <i>n</i>	Set the starting count for <code>lnk#.dat</code> files produced by <code>-crawl</code> . Set to 0 by default.
-localhost	Disable URLs that point to remote hosts. Set to <code>off</code> by default.
-locexec	Enable local program execution from local files only. Set to <code>off</code> by default.
-mime_header	Include MIME headers and force source dump.
-minimal	Use minimal instead of valid comment parsing. Set to <code>off</code> by default.
-nested_tables	Use nested-tables logic. Set to <code>off</code> by default.
-newschunksiz= <i>n</i>	Set the number of articles in chunked news listings.
-newsmaxchunk= <i>n</i>	Set the maximum number of news articles in listings before chunking.
-nobold	Disable bold video attribute.
-nobrowse	Disable directory browsing.
-nocc	Disable Cc: prompts for self-copies of mailings. Set to <code>off</code> by default.
-nocolor	Disable color support.
-noexec	Disable local program execution. Set to <code>on</code> by default.
-nofilereferer	Disable transmission of <code>Referer</code> headers for file URLs. Set to <code>on</code> by default.
-nolist	Disable the link list feature in dumps. Set to <code>off</code> by default.
-nolog	Disable mailing of error messages to document owners. Set to <code>on</code> by default.
-nonrestarting_sigwinch	Make window size change handler non-restarting. Set to <code>off</code> by default.

-nopause	Disable forced pauses for statusline messages.
-noprint	Disable some print functions, such as <code>restrictions=print</code> . Set to <code>off</code> by default.
-noredir	Don't follow <code>Location:</code> redirection. Set to <code>off</code> by default.
-noreferer	Disable transmission of <code>Referer</code> headers. Set to <code>off</code> by default.
-noreverse	Disable reverse video attribute.
-nostatus	Disable the miscellaneous information messages. Set to <code>off</code> by default.
-nounderline	Disable underline video-attribute.
-number_fields	Force numbering of links as well as form input fields. Set to <code>off</code> by default.
-number_links	Force numbering of links. Set to <code>off</code> by default.
-partial	Display partial pages while downloading. Set to <code>on</code> by default.
-partial_thres= <i>n</i>	Set the number of lines to render before repainting display with partial-display logic. Set to <code>-1</code> by default, which disables this feature.
-pauth=id:pw	Set the authentication information for a protected proxy server.
-popup	Handle single-choice <code>SELECT</code> options via popup windows instead of lists of radio buttons. Set to <code>off</code> by default.
-post_data	Read data for post forms from <code>STDIN</code> , terminated by <code>---</code> .
-preparsed	Show parsed text/html MIME type with <code>-source</code> and in source view to visualize how Lynx behaves with invalid HTML. Set to <code>off</code> by default.
-prettysrc	Use syntax highlighting and hyperlink handling in source view. Set to <code>off</code> by default.

<code>-print</code>	Enable print functions, the opposite of <code>-noprint</code> . Set to <code>on</code> by default.
<code>-pseudo_inlines</code>	Use pseudo-ALTs for inlines with no ALT string. Set to <code>on</code> by default.
<code>-raw</code>	Use default setting of 8-bit character translations or CJK mode for the startup character set. Set to <code>off</code> by default.
<code>-realm</code>	Restrict access to URLs in the starting realm. Set to <code>off</code> by default.
<code>-reload</code>	Flush the cache on a proxy server (only the first document affected). Set to <code>off</code> by default.
<code>-restrictions=options</code>	Set restriction options. Use <code>-restrictions</code> with no parameters to see list.
<code>-resubmit_posts</code>	Force resubmissions (no-cache) of forms with method <code>POST</code> when the documents they returned are sought with the <code>PREV_DOC</code> command or from the History list. Set to <code>off</code> by default.
<code>-rlogin</code>	Disable <code>rlogin</code> feature. Set to <code>off</code> by default.
<code>-selective</code>	Require <code>.www_browsable</code> files to browse directories.
<code>-short_url</code>	Enable examination of beginning and end of long URL in status line. Set to <code>off</code> by default.
<code>-show_cursor</code>	When set to <code>off</code> , hide the cursor in the lower-right corner, otherwise show cursor. Set to <code>on</code> by default.
<code>-show_rate</code>	Display the transfer rate. Set to <code>on</code> by default.
<code>-soft_dquotes</code>	Use emulation of old Netscape and Mosaic bug, which treated <code>></code> as a co-terminator for double quotation marks and tags. Set to <code>off</code> by default.
<code>-source</code>	Dump the source of the first URL to STDOUT and exit.
<code>-stack_dump</code>	Disable SIGINT cleanup handler. Set to <code>off</code> by default.

<code>-startfile_ok</code>	Allow non-HTTP start page and home page with <code>-validate</code> . Set to <code>off</code> by default.
<code>-stdin</code>	Read startfile from STDIN. Set to <code>off</code> by default.
<code>-tagsoup</code>	Use TagSoup rather than SortaSGML parser. Set to <code>off</code> by default.
<code>-telnet</code>	Disable telnet sessions. Set to <code>off</code> by default.
<code>-term=term</code>	Specify the terminal type to emulate.
<code>-tlog</code>	Use a Lynx Trace Log for the current session. Set to <code>on</code> by default.
<code>-tna</code>	Use Textfields Need Activation mode. Set to <code>off</code> by default.
<code>-trace</code>	Use Lynx trace mode. Set to <code>off</code> by default.
<code>-trace_mask</code>	Customize Lynx trace mode. Set to 0 by default.
<code>-traversal</code>	Traverse all HTTP links derived from startfile.
<code>-trim_input_fields</code>	Trim input text/textarea fields in forms. Set to <code>off</code> by default.
<code>-underline_links</code>	Use an underline/bold attribute for links. Set to <code>off</code> by default.
<code>-underscore</code>	Use an underline format in dumps. Set to <code>off</code> by default.
<code>-use_mouse</code>	Enable mouse support. Set to <code>off</code> by default.
<code>-useragent=Name</code>	Set alternate Lynx User-Agent header.
<code>-validate</code>	Accept only <code>http</code> URLs (meant for validation) implies more restrictions than <code>-anonymous</code> , but redirection is allowed for <code>http</code> and <code>https</code> . Set to <code>off</code> by default.
<code>-verbose</code>	Use <code>[LINK]</code> , <code>[IMAGE]</code> , and <code>[INLINE]</code> comments with file names of these images. Set to <code>on</code> by default.
<code>-version</code>	Display Lynx version information.
<code>-vikeys</code>	Enable vi-like key movement. Set to <code>off</code> by default.

<code>-width=n</code>	Set the screen width for formatting of dumps. The default is 80 columns.
<code>-with_backspaces</code>	Emit backspaces in output if using <code>-dump</code> or <code>-crawl</code> parameters. Set to <code>off</code> by default.

As you can see, you can control just about any type of HTTP or HTML setting directly from the command line. For example, if you want to post data to a Web form using the HTTP POST method, you just include your data in the `-post-data` parameter. If you want to store the cookies received by the Website in a special location, you use the `-cookie_save_file` parameter.

Many of the command line parameters define behaviors that control Lynx when you're using it in full-screen mode, allowing you to customize the behavior of Lynx as you're traversing Web pages.

There are often groups of command line parameters that you find useful in your normal browsing environment. Instead of having to enter these parameters on the command line every time you use Lynx, Lynx provides a general configuration file that defines the base behavior when you use Lynx. This configuration file is discussed in the next section.

The Lynx Configuration File

The `lynx` command reads a configuration file for many of its parameter settings. By default, this file is located at `/usr/local/lib/lynx.cfg`, although you'll find that many Linux

distributions change this to the `/etc` directory (`/etc/lynx.cfg`) (the Ubuntu distribution places the `lynx.cfg` file in the `/etc/lynx-curl` folder).

The `lynx.cfg` configuration file groups related parameters into sections to make finding parameters easier. The format of an entry in the configuration file is:

```
PARAMETER: value
```

where `PARAMETER` is the full name of the parameter (often, but not always in uppercase letters), and `value` is the value associated with the parameter.

Perusing this file, you'll find many parameters that are similar to the command line parameters, such as the `ACCEPT_ALL_COOKIES` parameter, which is equivalent to setting the `-accept_all_cookies` command line parameter.

There are also a few configuration parameters that are similar in function but different in name. The `FORCE_SSL_COOKIES_SECURE` configuration file parameter setting can be overridden by the `-force_secure` command line parameter.

However, you'll also find quite a few configuration parameters that don't match with command line parameters. These values can only be set from the configuration file.

The most common configuration parameters that you can't set on the command line are for the *proxy servers*. Some networks (especially corporate networks) use a proxy server as a middleman between the client's browser and the destination Website server. Instead of sending

HTTP requests directly to the remote Web server, client browsers must send their requests to the proxy server. The proxy server in turn sends the requests to the remote Web server, retrieves the results, and forwards them back to the client browser.

While this may seem like a waste of time, it's a vital function in protecting clients from dangers on the Internet. A proxy server can filter inappropriate content and malicious coding, or even detect sites used for Internet data phishing schemes (rogue servers pretending to be someone else in order to capture customer data). Proxy servers can also help reduce Internet bandwidth usage as they cache commonly viewed Web pages and return them to clients instead of having to download the original page again.

The configuration parameters used to define proxy servers are:

```
http_proxy:http://some.server.dom:port/  
https_proxy:http://some.server.dom:port/  
ftp_proxy:http://some.server.dom:port/  
gopher_proxy:http://some.server.dom:port/  
news_proxy:http://some.server.dom:port/  
newspost_proxy:http://some.server.dom:port/  
newsreply_proxy:http://some.server.dom:port/  
snews_proxy:http://some.server.dom:port/  
snewspost_proxy:http://some.server.dom:port/  
snewsreply_proxy:http://some.server.dom:port/  
nntp_proxy:http://some.server.dom:port/  
wais_proxy:http://some.server.dom:port/  
finger_proxy:http://some.server.dom:port/  
cso_proxy:http://some.server.dom:port/  
no_proxy:host.domain.dom
```

You can define a different proxy server for any network protocol supported by Lynx. The `NO_PROXY` parameter is a comma-separated list of Websites that you prefer to have

direct access to without using the proxy server. These are often internal Websites that don't require filtering.

The Lynx Environment Variables

As you can see from the wealth of command line options and configuration file parameters, the Lynx program is extremely customizable. However, the customization doesn't stop there. You can override many configuration file parameters using environment variables. If you're working in an environment where you don't have access to the `lynx.cfg` configuration file, you can override some default parameters by setting your local environment variables. [Table 24.2](#) lists the more common Lynx environment variables that you might need to use when using Lynx in a restricted environment.

[Table 24.2](#) The Lynx Environment Variables

Variable	Description
LYNX_CFG	Specify the location of an alternate configuration file.
LYNX_LSS	Specify the location of the default Lynx character set style sheet.
LYNX_SAVE_SPACE	Specify the location for saving files to disk.
NNTPSERVER	Specify the server to use for retrieving and posting USENET news.
PROTOCOL_PROXY	Override the proxy server for the specified protocol.
SSL_CERT_DIR	Specify the directory containing trusted certificates for accessing trusted sites.
SSL_CERT_FILE	Specify the file containing your trusted certificates.

You set these environment variables just as you would any other environment variable before using the Lynx program:

```
$ http_proxy=http://myproxy.com:8080
$ lynx
```

To specify a proxy server, you must provide the protocol, the server name, and the port used to communicate with the proxy server. If you need this variable setting, it's usually a good idea to include it in a common startup file for your shell (such as the `.bashrc` file for the bash shell) so you don't have to enter it every time.

Capturing Data from Lynx

When you use Lynx in a shell script, most likely you're trying to obtain a specific piece (or pieces) of information from a Web page. The technique to accomplish this is called *screen scraping*. In screen scraping, you're trying to programmatically find data in a specific location on a graphical screen so you can capture it and use it in your shell script.

The easiest way to perform screen scraping with `lynx` is to use the `-dump` option. This option doesn't bother trying to display the Web page on the terminal screen. Instead, it displays the Web page text data directly to `STDOUT`:

```
$ lynx -dump http://localhost/RecipeCenter/
The Recipe Center
    "Just like mom used to make"
Welcome
```



```
[1]Home
[2]Login to post
[3]Register for free login
```

```
[4]Post a new recipe
```

Each link is identified by a tag number, and Lynx displays a listing of all the tag references after the Web page data.

Once you have all of the text data from the Web page, you probably know what tools we're going to get out of the toolbox to start work on extracting data. That's right, our old friends the `sed` and `gawk` programs (see Chapter 18).

First, let's find some interesting data to collect. The Yahoo! weather Web page is a great source for finding the current weather conditions anywhere in the world. Each location uses a separate URL to display weather information for that city (you can find the specific URL for your city by going to the site in a normal browser and entering your city's information). The `lynx` command for finding the weather in Chicago, Illinois, is:

```
lynx -dump http://weather.yahoo.com/united-states/illinois/chicago-2379574/
```

This command dumps lots and lots of data from the Web page. The first step is to find the precise information you want. To do that, redirect the output from the `lynx` command to a file, and then search the file for your data. After doing that with the preceding command, we found this text in the output file:

```
Current conditions as of 1:54 pm EDT
Mostly Cloudy
```

```
Feels Like:
32 °F
```

```
Barometer:
    30.13 in and rising

Humidity:
    50%

Visibility:
    10 mi

Dewpoint:
    15 °F

Wind:
    W 10 mph
```

That's just about all the information about the current weather that you'll need. There's just one small problem with this output. You'll notice that the numbers are on a line below the heading. Trying to just extract individual numbers will be difficult. Chapter 18 discussed how to deal with a problem just like this.

The key to solving this is to write a sed script that can search for the data heading first. When you find it, you can then go to the correct line to extract the data. We're fortunate in this example in that all of the data we need are on lines by themselves. We should be able to solve this with just the sed script. If there had also been other text on the same line, we'd need to get out the gawk tool to filter out just the data we needed.

First, you need to create a sed script that will look for the `Current conditions` text, and then skip to the next line to get the text that describes the current weather condition, and then print it. Here's what that looks like:

```
$ cat sedcond
/Current conditions/{
n
p
}
```

```
}  
$
```

The address specifies to look for the line with the desired text. If the `sed` command finds it, the `n` command skips to the next line, and then the `p` command prints the contents of the line, which is the text describing the current weather conditions of the city.

Next, you'll need a `sed` script that can search for the `Feels Like:` text and then go to the next line to print the temperature:

```
$ cat sedtemp  
/Feels Like:/{  
n  
p  
}  
$
```

Perfect. Now, you can use these two `sed` scripts in a shell script that first captures the `lynx` output of the Web page to a temporary file, and then applies the two `sed` scripts to the Web page data to extract only the data you're looking for. Here's an example of how to do that:

```
$ cat weather  
#!/bin/bash  
# extract the current weather for Chicago, IL
```

```
URL="http://weather.yahoo.com/united-states/illinois/chicago-2379574/"  
LYNX='which lynx'  
TMPFILE='mktemp tmpXXXXXX'  
$LYNX -dump $URL > $TMPFILE  
conditions='cat $TMPFILE | sed -n -f sedcond'  
temp='cat $TMPFILE | sed -n -f sedtemp'  
rm -f $TMPFILE  
echo "Current conditions: $conditions"  
echo The current temp outside is: $temp  
$ ./weather  
Current conditions: Mostly Cloudy  
The current temp outside is: 32 °F  
$
```

The weather script connects to the Yahoo! weather Web page for the desired city, saves the Web page to a temporary file, extracts the appropriate text, removes the temporary file, and then displays the weather information. The beauty of this is that once you've extracted the data from a Website, you can do whatever you want with it, such as create a table of temperatures. You can then create a cron job (see Chapter 15) that runs every day to track daily temperatures.

Caution

The Internet is a dynamic place. Don't be surprised if you spend hours working out the precise location of data on a Web page, only to find that it's moved a couple of weeks later, breaking your scripts. In fact, it's quite possible that this example won't work by the time you read this book. The important thing is to know the process for extracting data from Web pages. You can then apply that principle to any situation.

The cURL Program

The popularity of Lynx has spawned another similar product, called cURL. The cURL program allows you to automatically transfer files from the command line using a specified URL. It currently supports the FTP, FTPS, HTTP, HTTPS, SCP, SFTP, TFTP, telnet, DICT, LDAP, LDAPS, and FILE protocols as specified in URLs.

While cURL isn't used as a Web page browser per se, it allows you to easily send or retrieve data unattended directly from the command line, or your shell scripts, just by using a simple command. This provides a great tool to have in your shell scripting toolbox.

This section walks you through the process of installing and using cURL in your shell scripts.

Caution

There's also a programming language called `curl`, owned and marketed by the Sumisho Computer System Corporation. Don't confuse cURL with the `curl` programming language.

Installing cURL

With its growing popularity, cURL is installed by default in many Linux distributions. With Ubuntu, you have to manually install it from the software repository. Just use this command:

```
$ sudo apt-get install curl
```

If cURL is not available for your Linux distribution, or if you just want to use the latest version, you can download the source code from the curl.haxx.se Website and compile it on your Linux system.

Again, the standard disclaimers apply; you must have the C development libraries installed on your Linux system for

this to work.

The next section describes how to use the curl program from the command line.

Exploring with cURL

By default, cURL returns the complete HTML code for the Web page to `STDOUT`:

```
$ curl http://www.google.com
<!doctype html><html><head><meta http-equiv="content-type" content=
"text/html; charset=ISO-8859-1"><title>Google</title>
[ listing truncated ]
```

As with the Lynx program, you can use the standard shell scripting techniques to extract individual data elements from the dumped Web page.

We like to use cURL for batch downloading files. These days, it seems like we're constantly downloading the latest Linux distribution ISO file. Because the ISO files are so large, we need to start the download and walk away while it's downloading. Once we know the URL for an ISO file, we can create a simple shell script using cURL to automate the process:

```
$ cat downld
#!/bin/bash
# download latest cURL file automatically
curl -s -o /home/rich/curl-7.18.0.tar.gz
http://curl.haxx.se/download/curl-7.18.0.tar.gz
$
```

The `-s` command line option puts curl in silent mode, not sending any data to `STDOUT`. The `-o` command line option redirects the output to a file name. While you can't tell from the listing, the `curl` command is all one line in the script. This

simple script just goes out and directly downloads the file from the cURL Website. You can now use the `at` or `cron` commands (see Chapter 15) to schedule the download to take place in the evening when you're not using your PC or network.

Networking with zsh

Chapter 22 describes all of the features available in the zsh shell. The zsh shell is a newer shell available for the Linux and Unix environments. One of the features of the zsh shell is plug-in modules. Instead of trying to combine lots of features into the core zsh shell, the zsh shell uses specialty modules, so you can pick and choose what commands you need to load. One of those modules is the TCP module.

The TCP module in the zsh shell provides some pretty amazing network capabilities directly from the command line. You can create a full TCP network session with another network device directly from your command line (or shell script). This section discusses the features of the zsh TCP module, and shows a simple client/server application you can build using a zsh shell script.

The TCP Module

The zsh shell uses modules to add additional features to the core zsh shell. Each module contains built-in commands that specialize in a specific area. The TCP

module provides built-in commands for a plethora of networking features.

To install the TCP module in a zsh shell, do the following:

```
% zmodload zsh/net/tcp
%
```

And that's all there is to install the module libraries in the shell! If you use the TCP module in your shell script, remember to include this line in your script. The module applies only to the current shell.

Once you load the TCP module, you have access to the `ztcp` command. The format for the `ztcp` command is:

```
ztcp [-acflLtv] [ -d fd] [args]
```

The available command line options are:

- `-a`: Accept a new connection.
- `-c`: Close an existing connection.
- `-d`: Use the specified file descriptor for the connection.
- `-f`: Force a connection to close.
- `-l`: Open a new socket for listening.
- `-L`: List currently connected sockets.
- `-t`: Exit if no connection is pending.
- `-v`: Display verbose information about the connection.

The `ztcp` command uses a file descriptor to interact with an open TCP connection. By default, the zsh referenced the file descriptor using the environment variable `$RESULT`. All you need to do is send data to the file descriptor specified in the `$RESULT` variable, and the TCP module will forward it to the remote host. Likewise, if the remote host sends you any

data, all you need to do is read it from the file descriptor specified in the `$RESULT` variable. Network programming doesn't get any easier than that!

The Client/Server Paradigm

Before diving into creating a client/server program using the `zsh` shell, it's a good idea to have an understanding of how exactly client and server programs operate. Obviously, they each have different responsibilities in the connection and transfer of data.

A *server* program listens to the network for requests coming from *clients*. A client initiates a request to the server for a connection. Once the server accepts the connection request, a two-way communication channel is available for each device to send and receive data. This process is shown in [Figure 24.2](#).

[Figure 24.2](#) The client/server communication diagram

Server**Client**

Listen on TCP port

Accept connection

send data

receive data

close connection

Connect to hostname and TCP port

receive data

send data

close connection

As you can see from [Figure 24.2](#), the server must perform two functions before it can communicate with the client. First, it must set up a specific TCP port to listen for incoming requests. When a connection request comes in, it must accept the connection.

The client's responsibility is much simpler. All it must do is attempt to connect to a server on the specific TCP port on which the server is listening. If the server accepts the connection, the two-way communication is available and data can be sent.

Once a connection is established between the server and the client, there must be some sort of communication process (or rule) used between the two devices. If both

devices attempt to listen for a message at the same time, they'll deadlock and nothing will happen. Likewise, if they both attempt to send a message at the same time, nothing will be accomplished.

It's your job as the network programmer to decide the protocol rules that your client and server programs must follow.

Client/Server Programming with zsh

To demonstrate the creation of a client/server program using `ztcp`, let's set up a simple network application. The server program we'll create will listen for connection requests on TCP port 5150. When a connection request comes in, the server will accept it and then send a welcome message to the client.

The server program will then wait to receive a message from the client. If it receives a message, the server will display the message and then send the same message back to the client. After sending the message, the server will loop back to listen for another message. This loop will continue until the server receives a message that consists of the text `exit`. When that happens, the server will terminate the session.

The client program we'll create will send a connection request to the server on TCP port 5150. When a connection is established, the client will need to receive the server's welcome message.

After receiving the message, the client displays it, and then queries the user for data to send to the server. After getting the message from the user, the client program will send it to the server and wait to receive the message back. If the message comes back, the client displays the message and loops back to request another message from the user. This loop will continue until the user enters the text `exit`. When this occurs, the client sends the exit text to the server and then terminates the session.

The next sections show the server and client programs.

The Server Program

Here's the code for the server program:

```
% cat server
#!/bin/zsh
# zsh TCP server script
zmodload zsh/net/tcp
ztcp -l 5150
fd=$REPLY

echo "Waiting for a client..."
ztcp -a $fd
clientfd=$REPLY
echo "client connected"

echo "Welcome to my server" >& $clientfd

while [ 1 ]
do
    read line <& $clientfd
    if [[ $line = "exit" ]]
    then
        break
    else
        echo Received: $line
        echo $line >& $clientfd
    fi
done
echo "Client disconnected session"
```

```
ztcp -c $fd
ztcp -c $clientfd
%
```

The server program follows the client/server paradigm shown in [Figure 24.2](#). It first uses the `-l` parameter to listen on the specified port (5150). The `$RESULT` variable contains the file descriptor that the Linux system returns to identify the connection. The server uses the `-a` parameter to accept a new connection request. This command waits until a new connection request comes in (called *blocking*). The script won't progress unless a connection request is accepted.

Each client connection uses a separate file descriptor from the listening port file descriptor. This allows you to maintain multiple client connections if you so desire (we don't in this simple exercise).

After accepting the connection, the server redirects a welcome message to the client's file descriptor:

```
echo "Welcome to my server" >& $clientfd
```

The zsh shell TCP module handles all of the mechanics of ensuring that the data is sent to the remote client.

Next, the server program enters an endless loop. It uses the `read` command to wait for data to come back from the client:

```
read line <& $clientfd
```

This command also blocks execution of the script until it receives data from the client. This could be a bad thing if the client has lost connection to the network. To prevent this problem, you can use the `-t` option in the `read` line to specify a timeout value (in seconds). If the server doesn't

receive data from the client in the timeout period, it continues on.

If the server receives data from the client, it displays the data on `STDOUT` and then sends it back to the client. If the data is equal to the text string `exit`, the server exits the loop, and uses the `-c` parameter in `ztcp` to close both the client's file descriptor and the listening port's file descriptor. If you prefer to have your server listen for another connection, after you close the client's file descriptor you can loop back to waiting to accept a new connection.

The Client Program

Here's the code for the client shell script program:

```
% cat client
#!/bin/zsh
# zsh TCP client program
zmodload zsh/net/tcp

ztcp localhost 5150
hostfd=$REPLY

read line <& $hostfd
echo $line

while [ 1 ]
do
    echo -n "Enter text: "
    read phrase
    echo Sending $phrase to remote host...
    echo $phrase >& $hostfd
    if [[ $phrase = "exit" ]]
    then
        break
    fi
    read line <& $hostfd
    echo "    Received: $line"
done
ztcp -c $hostfd
%
```

The client program must specify the IP address (or hostname) of the system the server program is running on and the proper TCP port number the server is listening to. When the server accepts the connection, the `ztcp` program sets the file descriptor for the connection and saves the value in the `$REPLY` variable. The client program reads the server's welcome message and then displays it:

```
read line <& $hostfd
echo $line
```

Next, the client program enters a `while` loop, querying the user for text to send to the server, reading the text entered, and sending the text to the server. After sending the text, it checks to see if the text entered was `exit`. If so, it breaks out of the loop and closes the file descriptor, which closes the TCP connection. If the text wasn't `exit`, it waits for the response from the server and then displays it.

Running the Programs

You can either run these programs on two separate Linux systems on your network or from two different terminal sessions on the same system. You must start the server program first so that it's available to listen for incoming connections when the client starts:

```
% ./server
Waiting for a client...
```

Then you can start the client:

```
% ./client
Welcome to my server
Enter text: test
Sending test to remote host...
```

```
Received: test
```

When the client connects, you'll see this on the server:

```
client connected  
Received: test
```

This will continue until the user enters the text `exit` on the client:

```
Enter text: exit  
Sending exit to remote host...  
%
```

Then you should see the server automatically exit:

```
Client disconnected session  
%
```

You now have the beginnings of a full-fledged network program! With the `zsh` shell and the `TCP` module, it's easy to send data between shell scripts operating on separate systems on a network.

Summary

This chapter walked you through the world of interfacing shell scripts to the Internet. One of the most popular tools for doing that is the `Lynx` program. `Lynx` is a command line program that can display Website information in a terminal session using text-mode graphics. In addition to that feature, `Lynx` also provides a way to retrieve just the raw data from a Website, and display it to `STDOUT`. You can use `Lynx` to extract data from a Website and then parse the data using standard Linux text-handling tools such as `sed` and `gawk` to find specific information.

The cURL program is another handy tool to have when interfacing with the Internet. The cURL program also allows you to dump data from a Website, and it provides a way to easily script file downloads from many different types of servers.

Finally, the chapter showed you how to use the TCP module of the zsh shell to write your own network programs. The zsh shell provides an easy way to communicate between shell scripts that are located on separate systems on the network.

In the next chapter, we examine how to utilize e-mail in your shell scripts. Often when you use shell scripts to automate processes, it's nice to get a message indicating whether the process failed or succeeded. By knowing how to work with the installed e-mail software on your system, you can easily send automated messages to anyone in the world.

Chapter 25

Using E-mail

In This Chapter

- E-mail and Linux
- Setting up an e-mail server
- Sending simple messages
- Using attachments with Mutt

With the popularity of e-mail, these days just about everyone has an e-mail address. Because of that, people often expect to receive data via e-mail instead of seeing files or printouts. That's no different in the shell scripting world. If you generate any type of report from your shell script, most likely at some point you'll be asked to e-mail the results to someone. This chapter shows you just how to set up your Linux system to support e-mailing directly from your shell scripts. It also shows you how to make sure that your Linux system can send outbound mail messages and how to make sure you have a mail client that can do that from the command line. But first, the chapter presents a brief overview of the way Linux handles e-mails in general.

The Basics of Linux E-Mail

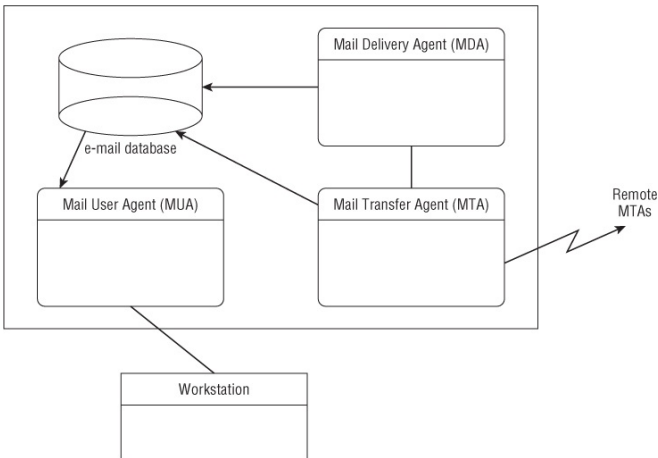
Sometimes the hardest part of using e-mail in your shell scripts is understanding just how the e-mail system works in Linux. Knowing what software packages perform particular tasks is crucial for getting e-mails from your shell script to your Inbox. This section walks you through the basics of how Linux systems use e-mail and what you need to have in place before you can use it.

E-Mail in Linux

One of the main goals of the Unix operating system was to modularize software. Instead of having one monolithic program that handles all of the required pieces of a function, Unix developers created smaller programs that each handle a smaller piece of the total functionality of the system.

This philosophy was used when implementing the e-mail systems used in Unix and was carried over to the Linux environment. In Linux, e-mail functions are divided into separate pieces, each assigned to separate programs. [Figure 25.1](#) shows how most open source e-mail software modularizes e-mail functions in the Linux environment.

[Figure 25.1](#) Linux modular e-mail environment



As you can see in [Figure 25.1](#), in the Linux environment the e-mail process is normally divided into three functions:

- The Mail Transfer Agent (MTA)
- The Mail Delivery Agent (MDA)
- The Mail User Agent (MUA)

The lines between these three functions are often fuzzy. Some e-mail packages combine functionality for the MDA and MTA functions, whereas others combine the MDA and MUA functions. The following sections describe these basic e-mail components and how they are implemented in Linux systems in more detail.

The Mail Transfer Agent

The MTA software is the core of the Linux e-mail system. It's responsible for handling both incoming and outgoing mail messages in the system. For each outgoing mail message the MTA must determine the destination of the recipient addresses. If the destination host is the local system, the MTA can either deliver it to the local mailbox directly or pass the message off to the local MDA for delivery.

However, if the destination host is a remote mail server, the MTA must establish a communication link with the MTA software on the remote host to transfer the message. There are two common methods that MTA software packages use to deliver mail to remote hosts:

- Direct delivery
- Proxy delivery

If your Linux system is directly connected to the Internet, it can often deliver messages destined for recipients on remote hosts directly to the remote host. The MTA software uses the Domain Name System (DNS) to resolve the proper network IP address to deliver the mail message, and then establishes the network connection using the Simple Mail Transfer Protocol (SMTP)

There are plenty of times when a host is not directly connected to the Internet, or it doesn't want to communicate directly with other remote hosts. In those situations, it usually uses a *smart host*. The smart host is a proxy server

that accepts mail messages from your Linux system and then attempts to directly deliver them to the intended recipient.

Note

Smart hosts are becoming more difficult to work with on the Internet because of spam. A rogue server can send spam by bouncing thousands of unsolicited commercial e-mail (UCE) messages off of a smart host to hide its identity. Most smart hosts now require some type of authentication before forwarding messages to other hosts.

For incoming messages, the MTA must be able to accept connection requests from remote mail servers and receive messages destined for local users. Again, the most common protocol used for this process is SMTP.

The Linux environment has many different types of open source MTA programs. Each program offers different features that distinguish it from the others. By far the two most popular you'll run into are:

- sendmail
- Postfix

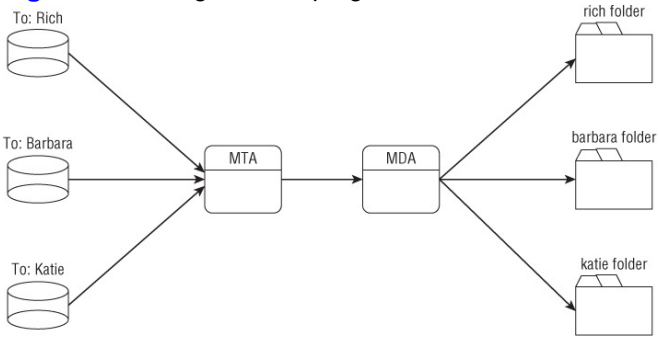
We'll examine both of these e-mail MTA packages in detail in the “Setting Up Your Server” section.

The Mail Delivery Agent

The MDA program's responsibility is to deliver a message

destined for a local user. It receives messages from the MTA program and must determine exactly how and where those messages should be delivered. [Figure 25.2](#) demonstrates how the MDA program interacts with the MTA program to deliver mail.

[Figure 25.2](#) Using an MDA program on a mail server



While sometimes the MDA function is performed within the MTA program itself, often Linux e-mail implementations rely on a separate standalone MDA program to deliver messages to local users. Because these MDA programs concentrate only on delivering mail to local users, they can add additional bells and whistles that aren't available on MTA programs that include MDA functionality. This enables the mail administrator to offer additional mail features to mail users, such as mail filtering for spam, out-of-office redirections, and automatic mail sorting.

When the MDA program receives a message, it must ensure that the message is delivered to the proper location, either to the local user's mailbox or to an alternate location defined by the local user.

There are currently three different types of user mailboxes commonly used on Linux systems:

- `/var/spool/mail` or `/var/mail` files
- `$HOME/mail` files
- Maildir-style mailbox directories

Each mailbox type has its own features that make it attractive to use. Most Linux distributions use either the `/var/spool/mail` or `/var/mail` directories to contain individual mailbox files, one file for each user account on the system. This is a central location for all mailbox files so that MUA programs know where to find everyone's mailbox file.

A few Linux distributions allow you to move the individual mailbox files to each user's `$HOME` directory. This provides greater security in that each mailbox file is located in an area already set with the proper access privileges.

Maildir-style mailboxes are a relatively new feature supported by some more advanced MTA, MDA, and MUA applications. Instead of each message being part of a mailbox file, the mailbox is a directory, and each message is a separate file in that directory. This helps cut down on mailbox corruption, as a single message won't corrupt the entire mailbox.

While Maildir-style mailbox directories offer increased performance, security, and fault tolerance, there are many

popular MDA and MUA programs that aren't able to use them. Just about all MDA and MUA programs can use the `/var/spool/mail` mailbox files.

Note

The original Unix location for mailboxes is `/var/spool/mail`. Most Linux distributions use this file-naming convention; however, there are a few Linux distributions that use `/var/mail` instead.

If your system does use a special MDA program to process incoming mail messages, most likely it's the popular Procmail program. Procmail allows each individual user to create a customized configuration file to define mail filters, out-of-office destinations, and separate mailboxes.

The Mail User Agent

So far we've followed the e-mail traffic from the remote host to the local host to an individual user's mailbox. The next step in the process is to allow individual users to view their e-mail messages.

The Linux e-mail model uses a local mailbox file or directory for each user to hold messages for that user. The job of the MUA program is to provide a method for users to interface with their mailboxes to read their messages.

It's important to remember that MUAs don't receive messages; they only display messages that are already in

the mailbox. Many MUA programs also offer the ability to create separate mail folders so the user can move mail from the default mailbox (often called the *Inbox*) to separate folders for organization.

Most MUA programs also provide the ability to send messages. This part gets a little fuzzy because, as you've already seen, sending e-mail messages is the job of the MTA program.

To perform this function, most MUA programs utilize the smart host feature in SMTP. Either the MUA program automatically delivers messages to the local MTA program for delivery, or you must define a remote smart host in the MUA configuration for it to send messages to for delivery.

Throughout the years, many different open source MUA programs have been available for the Linux platform. The following sections describe some of the more popular MUA programs you'll run across in Linux.

Mailx

The Mailx program is the most popular command line MUA program in use for the Linux environment. In all installations, the Mailx program installs with the executable file `mail`, indicating that it's a replacement for the mail program, rather than a separate program.

Note

The graphical desktop-oriented Linux distributions

such as Ubuntu and openSUSE don't install the command line Mailx program by default. You'll have to manually install it (see Chapter 8) as part of a mail client package. However, you have to be careful, because not all Linux mail packages are alike. For Ubuntu, the package is called mailutils, and uses the GNU Mailutils package instead of Mailx. It provides the same functionality, but with slightly different command line parameters.

The Mailx program allows users to access their mailboxes to read stored messages, as well as to send messages to other mail users, all from the command line. Here's a sample Mailx session.

```
$ mail
"/var/mail/rich": 2 messages 2 new
>N   1 Rich Blum           Thu Dec  9 10:07  13/579   Test message
   N   2 Rich Blum           Thu Dec  9 10:08  13/593   This is another test
? 1
Return-Path: <rich@rich-Parallels-Virtual-Platform>
X-Original-To: rich@rich-Parallels-Virtual-Platform
Delivered-To: rich@rich-Parallels-Virtual-Platform
Received: by rich-Parallels-Virtual-Platform (Postfix, from userid 1000)
        id 5C03F2606CF; Thu,  9 Dec 2010 10:07:48 -0500 (EST)
To: <rich@rich-Parallels-Virtual-Platform>
Subject: Test message
X-Mailer: mail (GNU Mailutils 2.1)
Message-Id: <20101209150748.5C03F2606CF@rich-Parallels-Virtual-Platform>
Date: Thu,  9 Dec 2010 10:07:48 -0500 (EST)
From: rich@rich-Parallels-Virtual-Platform (Rich Blum)

This is a test message
? d
? q
Held 1 message in /var/mail/rich
$
```

The first line shows the Mailx program being executed with no command line options. By default, this allows the user to check the messages in his mailbox. After entering

the `mail` command, a summary of all of the messages in the user's mailbox is displayed. The Mailx program can only read messages in the `/var/mail` format or `$HOME/mail` format. It's not able to process mail using the Maildir mail folder format.

Each user has a separate file that contains all of his messages. For some Linux distributions, the mailbox file doesn't exist until after the user account receives a message. The file name is usually the system login name of the user and is located in the system mailbox directory. Thus, all messages for username *rich* are stored in the file `/var/mail/rich` on the Linux system. As new messages are received for the user, they are appended to the end of the file.

You can also send mail messages using the mail command line program:

```
$ mail barbara
Cc:
Subject: This is a test sent to Barbara
Hello Barbara -
This is a test message I'm sending from the command line.
$
```

The name of the recipient is included on the command line with the program name. The Mailx program queries for any addresses to send copies to (the CC: prompt) and the message subject. It then allows you to type in the text of the message. To terminate the message, press Ctrl+D. The Mailx program then attempts to pass the message to the MTA program for delivery.

Mutt

As advancements were made to the Unix environment, MUA programs became fancier. One of the first attempts at graphics on Unix systems was the ncurses graphics library. Using ncurses, a program could manipulate the location of a cursor on the terminal screen and place characters almost anywhere on the terminal.

One MUA program that takes advantage of the ncurses library is the Mutt program. When you start Mutt, it paints a user-friendly menu on the terminal display, listing the messages similar to the output of the Mailx program. You can select a message and view it in the display, as shown in [Figure 25.3](#).

[Figure 25.3](#) The Mutt program

```
rich@rich-Parallels-Virtual-Platform: ~
File Edit View Search Terminal Help
i:Exit -:PrevPg <Space>:NextPg v:View Attachm. d:Del r:Reply j:Next ?:Help
Date: Thu, 9 Dec 2010 10:19:50 -0500 (EST)
From: Rich Blum <rich@rich-Parallels-Virtual-Platform>
To: rich@rich-Parallels-Virtual-Platform
Subject: Test message
X-Mailer: mail (GNU Mailutils 2.1)

This is a test message.

-N F- 1/1: Rich Blum Test message -- (all)
```

The Mutt program uses key combinations to perform standard functions, such as read a message and start a new message. Possibly the most useful feature for shell script programmers is the ability to send messages directly from the command line, without going into text-graphics mode. We'll be covering the Mutt program in much more detail later on in “The Mutt Program” section.

Graphical E-Mail Clients

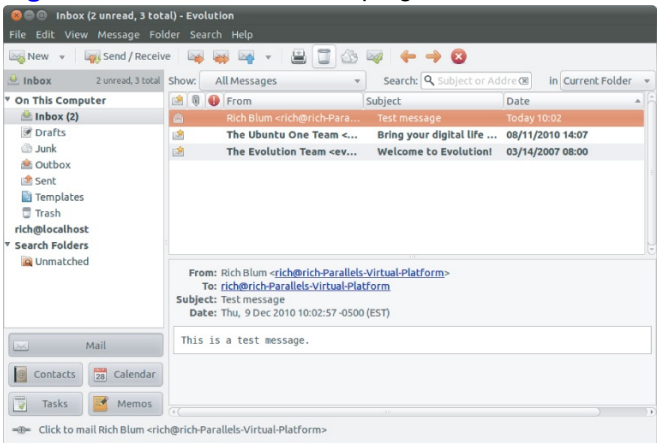
Almost all Linux systems support the graphical X Window environment. There are many e-mail MUA programs that utilize the X Window system to display message

information. The two most popular graphical MUA programs available are:

- KMail for the KDE windows environment
- Evolution for the GNOME windows environment

Each of these packages allows you to interact with your local Linux mailbox, as well as connect to remote mail servers to read mail messages. [Figure 25.4](#) shows a sample Evolution session screen.

Figure 25.4 The Evolution MUA program main screen



To connect with remote servers, both KMail and Evolution support the Post Office Protocol (POP) and the more advanced Internet Message Access Protocol (IMAP).

While the KMail and Evolution MUA programs are great for desktop Linux, they aren't so useful in shell scripting.

Setting Up Your Server

Before you can start sending your automated e-mail messages off to the universe, you'll need to ensure that your Linux system has an MTA package running and that it's configured correctly. This is no small task in itself, but fortunately, some Linux distributions provide some basic tools to help you out.

This section walks you through the basics of the two most popular e-mail MTA programs used in Linux: sendmail and Postfix. While there have been complete books written on properly configuring each of these packages, we'll just look at the basics to see how to get e-mails off the Linux system and into your Inbox.

sendmail

The sendmail MTA package is one of the most popular open source MTA packages used by Internet mail servers. In the past, it had been plagued with stories about backdoors and security flaws; however, it has been rewritten not only to remove the security flaws but also to incorporate many newer MTA features such as spam control. The newer versions of the sendmail program have proven to be secure as well as versatile.

Parts of the sendmail Program

The main executable program is called sendmail. It normally runs in background mode, listening for SMTP connections from remote mail servers and forwarding outbound messages from local users.

In addition to the main sendmail program, there's a configuration file and several tables that it uses to contain information used while processing incoming and outgoing mail messages. [Table 25.1](#) lists all the parts used in a normal sendmail installation.

Table 25.1 The sendmail Configuration Files

File	Description
sendmail.cf	Text file that controls the behavior of the sendmail program.
sendmail.cw	Text file that contains a list of domain names that the sendmail program will receive messages for.
sendmail.ct	Text file that contains a list of trusted users who can control the sendmail operations.
aliases	Binary file that contains a list of valid local mail addresses that can redirect mail to another user, a file, or a program.
newaliases	Executable program that creates a new aliases database file from a text file.
mailq	Executable program that checks the mail queue and prints any messages.
mqueue	A directory used to store messages waiting to be delivered.
mailertable	Text file used to specify route paths for specific domains.
domaintable	Text file used to map old domain names to new ones.
	Text file used to map users and domains to alternate

virtusertable	addresses.
relay-domains	Text file used to list specific hosts that are allowed to relay messages through the sendmail program.
access	Text file that lists specific domains from which received messages are either allowed or refused.

Unless you're running the main mail server for a corporation or Internet service provider, all you'll need to worry about is the `sendmail.cf` configuration file. In fact, many Linux distributions that use sendmail create and configure a core `sendmail.cf` configuration file for you automatically that should work just fine in most simple applications.

The sendmail.cf File

The sendmail program needs to be told how to handle messages as the server receives them. As an MTA, sendmail processes incoming mail and redirects it to another mail package, either on a remote system or on the local system. The configuration file is used to direct sendmail how to manipulate the destination mail addresses to determine where and how to forward the messages. The default location for the configuration file is `/etc/mail/sendmail.cf`.

The `sendmail.cf` file consists of rule sets that parse the incoming mail message and determine what actions to take. Each rule set is used to identify certain mail formats and instruct sendmail how to handle that message.

As the sendmail program receives a message, it parses the message header and passes the message through the

various rule sets to determine an action to take on the message. The sendmail configuration file includes rules that allow sendmail to handle mail in many different formats. Mail received from an SMTP host has different header fields than mail received from a local user. The sendmail program must know how to handle any mail situation.

Rules also have helper functions defined in the configuration file. There are three different types of helper functions that you can define:

- **Classes:** Define common phrases that are used to help the rule sets identify certain types of messages.
- **Macros:** Values that are set to simplify the typing of long strings in the configuration file.
- **Options:** Defined to set parameters for the sendmail program's operation.

The configuration file is made up of a series of classes, macros, options, and rule sets. Each function is defined as a single text line in the configuration file.

Each line in the configuration file begins with a single character that defines the action for that line. Lines that begin with a space or a tab are continuation lines from a previous action line. Lines that begin with a pound sign (#) indicate comments and are not processed by sendmail.

The action at the beginning of the text line defines what the line is used for. [Table 25.2](#) shows the standard sendmail actions and what they represent.

[Table 25.2](#) sendmail Configuration File Lines

Configuration Line	Description
C	Defines classes of text
D	Defines a macro
F	Defines files containing classes of text
H	Defines header fields and actions
K	Defines databases that contain text to search
M	Defines mailers
O	Defines sendmail options
P	Defines sendmail precedence values
R	Defines rule sets to parse addresses
S	Defines rule set groups

As we mentioned, most likely you won't have to start from scratch with your `sendmail.cf` configuration file; the Linux distribution should create a standard template file for you. Most likely, the only piece you'll have to worry about is if you must use a smart host to forward mail for you. The `DS` configuration line controls this feature:

DSmyisp.com

Just add the hostname of the smart host immediately after the `DS` tag.

Postfix

The Postfix software package is quickly becoming one of the more popular e-mail packages available for Unix and Linux systems. Postfix was developed by Wietse Venema to provide an alternate MTA for standard Unix-type servers.

The Postfix software is capable of turning any Unix or Linux system into a fully functional e-mail server.

It is the responsibility of the MTA package to manage messages that come into or leave the mail server. Postfix accomplishes this message tracking by using several different modular programs and a system of mail queue directories. Each program processes messages through the various message queues until they are delivered to their final destinations. If at any time the mail server crashes during a message transfer, Postfix can determine what queue the message was last successfully placed in and attempt to continue the message processing.

Parts of the Postfix System

The Postfix system consists of several mail queue directories and executable programs, all interacting with each other to provide mail service.

The Postfix package utilizes a master program that runs as a background process at all times. The master program allows Postfix to spawn programs that scan the mail queues for new messages and send them to the proper destinations.

The master program utilizes additional helper programs that are started as needed, depending on their function. They can be configured to remain running for set times after they are utilized. This allows the master program to reuse a running helper program if necessary, saving processing time. After a set time limit, the helper program quietly stops

itself.

The master program is used to control the overall operation of Postfix. [Table 25.3](#) shows the helper programs that Postfix uses to transfer mail messages.

Table 25.3 The Postfix Helper Programs

Program	Description
bounce	Posts a log in the bounce message queue for bounced messages and returns the bounced message to the sender
cleanup	Processes incoming mail headers and places messages in the incoming queue
error	Processes message delivery requests from qmgr, forcing messages to bounce
flush	Processes messages waiting to be retrieved by a remote mail server
local	Delivers messages destined for local users
pickup	Waits for messages in the maildrop queue and sends them to the cleanup program to begin processing
pipe	Forwards messages from the queue manager program to external programs
postdrop	Moves an incoming message to the maildrop queue when that queue is not writable by normal users
qmgr	Processes messages in the incoming queue, determining where and how they should be delivered, and spawns programs to deliver them
sendmail	Provides a sendmail compatible interface for programs to send messages to the maildrop queue
showq	Reports Postfix mail queue status
smtp	Forwards messages to external mail hosts using the SMTP protocol

smtpd	Receives messages from external mail hosts using the SMTP protocol
trivial-rewrite	Receives messages from the cleanup program to ensure header addresses are in a standard format for the qmgr program, and is used by the qmgr program to resolve remote host addresses

Postfix uses several different message queues for managing e-mail messages as they are processed. Each message queue contains messages in a different message state in the Postfix system. [Table 25.4](#) lists the message queues that are used by Postfix.

[Table 25.4](#) Postfix Message Queues

Queue	Description
maildrop	New messages waiting to be processed, received from local users.
incoming	New messages waiting to be processed or received from remote hosts, as well as processed messages from local users.
active	Messages that are ready to be delivered by Postfix.
deferred	Messages that have failed on an initial delivery attempt and are waiting for another attempt.
flush	Messages that are destined for remote hosts that will connect to the mail server to retrieve them.
mail	Delivered messages stored for local users to read.

If the Postfix system should be shut down at any time, messages remain in the last queue in which they were placed. When Postfix is restarted, it will automatically begin processing messages from the queues. This is a great feature of Postfix, making it one of the most robust e-mail

server systems available!

Postfix Configuration Files

A vital part of the Postfix system is the configuration files. Postfix uses three separate configuration files to allow you to set the parameters used to direct Postfix on how to handle messages. Unlike some other MTA programs, it's possible to change configuration information while the Postfix server is running and issue a command to have Postfix load the new information without completely stopping the mail server.

The three configuration files are normally stored in a common Postfix directory. Often, the default location for this directory is `/etc/postfix`. Usually, all users have access to view the configuration files, whereas only the root user has the ability to change values within the files. Of course, this can be modified for your own security situation. [Table 25.5](#) lists the Postfix configuration files.

[Table 25.5](#) Postfix Configuration Files

File	Description
<code>install.cf</code>	Contains information from the install parameters used when Postfix was installed.
<code>main.cf</code>	Contains parameters used by the Postfix programs when processing messages.
<code>master.cf</code>	Contains parameters used by the Postfix master program when running core programs.

The `install.cf` configuration file allows you to retrieve

installation parameters that were used when the Postfix software was first installed on the system. This provides an easy way to determine which features are or aren't available in the software setup.

The `master.cf` configuration file controls the behavior of the core Postfix programs. Each program is listed in a separate line along with the parameters to control its operation. Here's a sample `master.cf` file with default settings.

```
# =====
#service type private unpriv chroot wakeup maxproc command + args
#          (yes)    (yes)    (yes)    (never) (50)
# =====
smtp      inet  n       -       n       -       -       smtpd
pickup    fifo  n       -       n       60      1       pickup
cleanup   unix  n       -       n       -       0       cleanup
qmgr      fifo  n       -       n       300     1       qmgr
rewrite   unix  -       -       n       -       -       trivial-rewrite
bounce    unix  -       -       n       -       0       bounce
defer      unix  -       -       n       -       0       bounce
trac       unix  -       -       n       -       0       bounce
verify    unix  -       -       n       -       1       verify
flush     unix  n       -       n       1000    0       flush
proxymap  unix  -       -       n       -       -       proxymap
smtp       unix  -       -       n       -       -       smtp
relay      unix  -       -       n       -       -       smtp -o fallback_relay=
showq      unix  n       -       n       -       -       showq
error      unix  -       -       n       -       -       error
local      unix  -       n       n       -       -       local
virtual    unix  -       n       n       -       -       virtual
lmtp       unix  -       -       n       -       -       lmtp
anvil      unix  -       -       n       -       1       anvil
scache     unix  -       -       n       -       1       scache
```

The `master.cf` configuration file also includes lines for directing Postfix on how to interface with external MDA software, such as Procmail.

The Postfix operational parameters are set in the `main.cf` configuration file. All of the Postfix operational parameters

have default values that are implied within the Postfix system. If a parameter value is not present in the `main.cf` file, its value is preset by Postfix. If a parameter value is present in the `main.cf` file, its contents override the default value.

Each Postfix parameter is listed on a separate line in the configuration file along with its value, in the following form:

```
parameter = value
```

Both `parameter` and `value` are plain text strings that can be easily read and changed if necessary. The Postfix master program reads the parameter values in the `main.cf` file when Postfix is first started and again whenever a Postfix reload command is issued.

Two examples of Postfix parameters are the `myhostname` and `mydomain` parameters. If they are not specified in the `main.cf` configuration file, the `myhostname` parameter assumes the results of a `gethostname()` command on the Linux system, whereas `mydomain` assumes the domain part of the default `myhostname` parameter. Often a single mail server will handle mail for an entire domain. This is an easy setting in the Postfix configuration file:

```
myhostname = mailserver.smallorg.org  
mydomain = smallorg.org
```

When Postfix starts, it will recognize the local mail server as a mailserver.smallorg.org and the local domain as smallorg.org and will ignore any system set values.

If you need to specify a smart host, do that with the `relayhost` parameter:

```
relayhost = myisp.com
```

You can also specify an IP address here, but it must be enclosed in square brackets.

Sending a Message with Mailx

The main tool you have available for sending e-mail messages from your shell scripts is the Mailx program. Not only can you use it interactively to read and send messages, but you can also use the command line parameters to specify how to send a message.

The format for the Mailx program's command line for sending messages is:

```
mail [-eIinv] [-a header] [-b addr] [-c addr] [-s subj] to-addr
```

The `mail` command uses the command line parameters shown in [Table 25.6](#).

Table 25.6 The Mailx Command Line Parameters

Parameter	Description
-a	Specify additional SMTP header lines.
-b	Add a BCC: recipient to the message.
-c	Add a CC: recipient to the message.
-e	Don't send the message if it's empty.
-i	Ignore TTY interrupt signals.
-I	Force Mailx to run in interactive mode.
-n	Don't read the <code>/etc/mail.rc</code> startup file.
-s	Specify a Subject line.
-v	Display details of the delivery on the terminal.

As you can see from [Table 25.6](#), you can pretty much create an entire e-mail message just from the command line parameters. The only thing you need to add is the message body.

To do that, you need to redirect text to the `mail` command. Here's a simple example of how to create and send an e-mail message directly from the command line:

```
$ echo "This is a test message" | mail -s "Test message" rich
```

The Mailx program sends the text from the `echo` command as the message body. This provides an easy way for you to send messages from your shell scripts. Here's a quick example:

```
$ cat factmail
#!/bin/bash
# mailing the answer to a factorial

MAIL='which mail'

factorial=1
counter=1

read -p "Enter the number: " value
while [ $counter -le $value ]
do
    factorial=$((factorial * $counter))
    counter=$((counter + 1))
done

echo "The factorial of $value is $factorial" | mail -s "Factorial
answer" $USER
echo "The result has been mailed to you."
```

This script does not assume that the Mailx program is located in the standard location. It uses the `which` command to determine just where the `mail` program is.

After calculating the result of the factorial function, the

shell script uses the `mail` command to send the message to the user-defined `$USER` environment variable, which should be the person executing the script.

```
$ ./factmail
Enter the number: 5
The result has been mailed to you.
$
```

All you need to do is check your mail to see if the answer arrived:

```
$ mail
"/var/mail/rich": 1 message 1 new
>N      1 Rich Blum           Thu Dec  9 10:32  13/586   Factorial answer
?
Return-Path: <rich@rich-Parallels-Virtual-Platform>
X-Original-To: rich@rich-Parallels-Virtual-Platform
Delivered-To: rich@rich-Parallels-Virtual-Platform
Received: by rich-Parallels-Virtual-Platform (Postfix, from userid 1000)
         id B4A2A260081; Thu,  9 Dec 2010 10:32:24 -0500 (EST)
Subject: Factorial answer
To: <rich@rich-Parallels-Virtual-Platform>
X-Mailer: mail (GNU Mailutils 2.1)
Message-Id: <20101209153224.B4A2A260081@rich-Parallels-Virtual-Platform>
Date: Thu,  9 Dec 2010 10:32:24 -0500 (EST)
From: rich@rich-Parallels-Virtual-Platform (Rich Blum)

The factorial of 5 is 120
?
```

It's not always convenient to send just one line of text in the message body. Often, you'll need to send an entire output as the e-mail message. In those situations, you can always redirect text to a temporary file, and then use the `cat` command and redirect the output to the `mail` program.

Here's an example of sending a larger amount of data in an e-mail message:

```
$ cat diskmail
#!/bin/bash
# sending the current disk statistics in an e-mail message
```

```
date=`date +%m/%d/%Y`  
MAIL=`which mail`  
TEMP=`mktemp tmp.XXXXXX`
```

```
df -k > $TEMP  
cat $TEMP | $MAIL -s "Disk stats for $date" $1  
rm -f $TEMP
```

The `diskmail` program gets the current date using the `date` command (along with some special formatting), finds the location of the Mailx program, and then creates a temporary file. After all that, it uses the `df` command to display the current disk space statistics (see Chapter 4), redirecting the output to the temporary file.

It then redirects the temporary file to the `mail` command, using the first command line parameter for the destination address, and the current date in the Subject header. When you run the script, you won't see anything appear on the command line output:

```
$ ./diskmail rich
```

But if you check your mail, you should see the sent message:

```
$ mail  
"/var/mail/rich": 1 message 1 new  
>N 1 Rich Blum Thu Dec 9 10:35 19/1020 Disk stats for 12/09/2010  
?  
Return-Path: <rich@rich-Parallels-Virtual-Platform>  
X-Original-To: rich@rich-Parallels-Virtual-Platform  
Delivered-To: rich@rich-Parallels-Virtual-Platform  
Received: by rich-Parallels-Virtual-Platform (Postfix, from userid 1000)  
id 3671B260081; Thu, 9 Dec 2010 10:35:39 -0500 (EST)  
Subject: Disk stats for 12/09/2010  
To: <rich@rich-Parallels-Virtual-Platform>  
X-Mailer: mail (GNU Mailutils 2.1)  
Message-Id: <20101209153539.3671B260081@rich-Parallels-Virtual-Platform>  
Date: Thu, 9 Dec 2010 10:35:39 -0500 (EST)  
From: rich@rich-Parallels-Virtual-Platform (Rich Blum)
```

```
Filesystem 1K-blocks Used Available Use% Mounted on
```

/dev/sda1	63315876	2595552	57504044	5% /
none	507052	228	506824	1% /dev
none	512648	192	512456	1% /dev/shm
none	512648	100	512548	1% /var/run
none	512648	0	512648	0% /var/lock
none	4294967296	0	4294967296	0% /media/psf
?				

Now all you need to do is schedule the script to run every day using the cron feature, and you can get disk space reports automatically e-mailed to your inbox! System administration doesn't get much easier than that!

The Mutt Program

The Mutt program is another popular e-mail client package for the Linux command line, developed in 1995 by Michael Elkins. It has one feature that's not available in the Mailx program, which makes it a good tool to have handy for your shell scripts.

The Mutt program has the capability to send files as attachments in your e-mail messages. Instead of having to incorporate a long text file in the body of your e-mail message as we did with Mailx, you can use the Mutt program and include the text file as a separate attachment to the main message body. This feature is great for e-mailing long files, such as log files.

This section walks you through installing Mutt on your Linux system and using it to attach files to e-mail messages in your shell scripts.

Installing Mutt

The Mutt program is not a popular package in this day of fancy graphical e-mail clients such as KMail or Evolution, so it's a good bet that your Linux distribution doesn't have it installed by default. However, most Linux distributions include it in the normal distribution files for installation using the standard software installation methods (see Chapter 8). For the Ubuntu environment, you can install Mutt using the command line command:

```
sudo apt-get install mutt
```

If your Linux distribution doesn't include the Mutt package, or you just want to install the latest version, you can go to the Mutt Website (www.mutt.org), download the latest source code file, and then use the techniques discussed in Chapter 8 to compile and install the package from the source code package.

The Mutt Command Line

The `mutt` command provides parameters for you to use to control how Mutt operates. [Table 25.7](#) shows the command line parameters available to you.

Table 25.7 The Mutt Command Line Parameters

Parameter	Description
<code>-A alias</code>	Pass an expanded version of the specified alias to STDOUT.
<code>-a file</code>	Attach the specified file to your message using the MIME protocol.

-b <i>address</i>	Specify a blind-carbon-copy (BCC) recipient.
-c <i>address</i>	Specify a carbon-copy (CC) recipient.
-D	Print the value of all configuration options to STDOUT.
-e <i>command</i>	Specify a configuration command to be run after processing of initialization files.
-f <i>mailbox</i>	Specify a mailbox file to load.
-F <i>muttrc</i>	Specify an initialization file to read instead of \$HOME/.muttrc.
-h	Display help text.
-H <i>draft</i>	Specify a draft file that contains header and body to use to send a message.
-i <i>include</i>	Specify a file to include in the body of a message.
-m <i>type</i>	Specify a default mailbox type.
-n	Ignore the system configuration file.
-p	Resume a postponed message.
-Q <i>query</i>	Query a configuration variable value. The query is executed after all configuration files have been parsed and any commands given on the command line have been executed.
-R	Open a mailbox in read-only mode.
-s <i>subject</i>	Specify the subject of the message.
-v	Display the Mutt version number and compile-time definitions.
-x	Emulate the Mailx compose mode.
-y	Start with a listing of all mailboxes specified by the mailboxes command.
-z	When used with -f, don't start if there are no messages in the mailbox.
-Z	Open the first mailbox specified by the mailboxes command that contains new mail.

With the myriad of command line parameters you can customize your e-mail message directly from the command line, which is exactly what you'll want to do in your shell scripts.

Much as with the Mailx program, there's one thing that you can't specify on the command line with the Mutt program, and that's the message body text. If you don't redirect text to the Mutt program, it will start in text-graphics mode with an editor window for you to type the message body in.

This is not a good thing for the shell script, so you'll always want to redirect some type of text for the message body, even if you're using the attachment option to specify a file to attach. The next section demonstrates how to do this.

Using Mutt

Now you're ready to start using the Mutt program in your shell scripts. To create the basic `mutt` command in your shell script, you'll want to include command line options that specify the subject of the message, the attachment file, and all the recipients of the message:

```
mutt -s Subject -a file -- recipients
```

The `recipients` list is a space-separated list of the e-mail addresses to send the message to. If you want to attach more than one file, you can list them separated with spaces after the `-a` option; the `--` symbol is used to separate the file names from the list of recipient addresses. The `file`

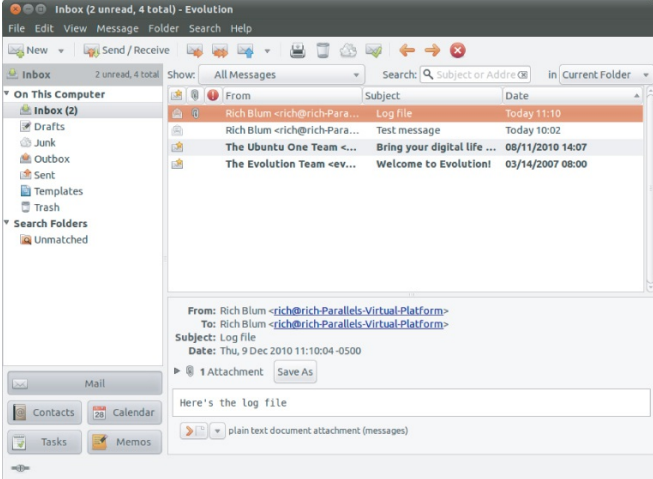
parameter must be an absolute pathname, or a relative pathname relative to the current working directory from which you're running the `mutt` command.

There's one other catch with the `mutt` command. If you don't redirect text for the message body, Mutt will automatically go to full-screen mode for you to enter the text in an editor window. Most likely this is not what you want to do, so be sure to redirect some text for the message body, even if it's an empty file:

```
# echo "Here's the log file" | mutt -s "Log file" -a  
/var/log/messages -- rich
```

This command sends the system log file as an attachment to the e-mail address `rich` on the local system. Note that you must also have the proper permissions to access the file you want to attach. [Figure 25.5](#) shows the received e-mail message in the Evolution mail client.

[Figure 25.5](#) Using Evolution to view a message with an attachment



Notice that the message includes the body text from the echo statement along with a separate icon for the attached file. You can save the attached file directly from the KMail client.

Caution

If you look at the name of the attached file, you'll notice that Mutt uses the basename of the attached file as the file name in the attachment. Be careful when using temporary files, as Mutt will use the temporary file name as the attachment file name. You'll be better off saving temporary files using more descriptive names rather than using temporary file names.

Summary

This chapter discusses how to incorporate e-mail in your shell scripts. The ability to regularly e-mail reports to customers is a great feature to offer in your shell scripts.

Before using the command line to send e-mail messages, you'll need to know how e-mail works in the Linux environment and what applications you'll need to have installed and configured. The Linux e-mail environment consists of three elements, a Mail Transport Agent (MTA) program, a Mail Delivery Agent (MDA) program (which often is part of the MTA program), and a Mail User Agent (MUA) program.

The MTA program is responsible for sending and receiving mail messages for the Linux system. It must know how to pass incoming mail messages to the proper user mailboxes, as well as how to send outbound mail messages destined for users on remote mail servers. Often, the MTA program will use a proxy server (also called a smart host) to do the detailed mail delivery. It forwards any message destined for a user on a remote mail server to the smart host for delivery. It can rely on the smart host to make the delivery on its own.

The MDA program is responsible for ensuring that mail destined for local users ends up in the proper local mailbox. Sometimes this function is performed directly by the MTA

program. However, if you require advanced mail delivery features, such as out-of-office notifications or spam filtering, you can configure the MTA program to pass messages to the MDA program, which often has these capabilities built in.

The MUA program allows individual system users to access messages in their mailboxes and pass outbound messages to other users to the MTA for delivery. These programs can range from simple command line programs, such as Mailx and Mutt, to fancy graphical programs such as KMail and Evolution.

The easiest way to send e-mails from your shell script is to use the Mailx program. This program allows you to specify the subject header and one or more recipients on the command line. You create the message body by redirecting text to the Mailx program. You can do this using either the `echo` command for single line text, or the `cat` command to redirect the contents of a file to the mail message.

The Mutt program is a more advanced command line MUA program that provides the ability to attach files to the mail message, rather than including text inside the message body. This allows you to attach large text files that your customers can easily save to disk for examination in other programs, such as spreadsheets or word processing packages.

The last two chapters in this book cover an important part of shell scripts, scripting administration functions. If you're a

Linux system administrator, most likely you'll run into a situation where you want to regularly monitor the status of a system feature. By creating a shell script and placing it in a cron job, you can easily monitor what's going on in your Linux system.

Chapter 26

Writing Script Utilities

In This Chapter

- Monitoring disk space
- Performing backups
- Managing user accounts

There is no place where shell script programming is more useful than writing script utilities for the Linux system administrator. The typical Linux system administrator has a myriad of jobs that need to be done daily, from monitoring disk space to backing up important files to managing user accounts. Shell script utilities can make the life of the system administrator much easier! This chapter demonstrates some of the capabilities you have writing script utilities in the bash shell.

Monitoring Disk Space

One of the biggest problems with multi-user Linux systems is the amount of available disk space. In some situations,

such as in a file sharing server, disk space can fill up almost immediately just because of one careless user.

This shell script utility helps you determine the top ten disk space consumers for designated directories. It produces a date-stamped report that allows disk space consumption trends to be monitored.

The Required Functions

The first tool you need to use is the `du` command (see Chapter 4). This command displays the disk usage for individual files and directories. The `-s` option lets you summarize totals at the directory level. This comes in handy when calculating the total disk space used by an individual user. Here's what it looks like to use the `du` command to summarize each user's `$HOME` directory for the `/home` directory contents:

```
$ du -s /home/*
6174428      /home/consultant
4740         /home/Development
4740         /home/Production
3860         /home/Samantha
7916         /home/Timothy
140376116    /home/user
$
```

The `-s` option works well for users' `$HOME` directories, but what if we wanted to view disk consumption in a system directory such as `/var/log`?

```
$ du -s /var/log/*
4          /var/log/alternatives.log
44         /var/log/alternatives.log.1
4          /var/log/apparmor
176        /var/log/apt
0          /var/log/aptitude
```

```

$ du -s /var/log/apititude.1.gz
4      /var/log/apititude.2.gz
4      /var/log/apititude.3.gz
4      /var/log/apititude.4.gz
4      /var/log/apititude.5.gz
160    /var/log/auth.log
...
$

```

The listing quickly becomes too detailed. The `-s` option works better for our purposes here, providing a total for each directory and subdirectory individually. This allows you to pinpoint problem areas quickly:

```

$ du -s /var/log
176    /var/log/apt
52     /var/log/exim4
1048   /var/log/dist-upgrade/20101011-1337
6148   /var/log/dist-upgrade
1248   /var/log/installer
228    /var/log/gdm
4      /var/log/news
4      /var/log/samba/cores/winbindd
4      /var/log/samba/cores
16     /var/log/samba
4      /var/log/unattended-upgrades
4      /var/log/sysstat
4      /var/log/speech-dispatcher
108    /var/log/ConsoleKit
64     /var/log/cups
4      /var/log/apparmor
12     /var/log/fsck
4844   /var/log
$

```

Because we are interested in the directories consuming the biggest chunks of disk space, you need to use the `sort` command (see Chapter 4) on the listing produced by `du`:

```

$ du -s /var/log | sort -rn
6148   /var/log/dist-upgrade
4864   /var/log
1248   /var/log/installer
1048   /var/log/dist-upgrade/20101011-1337
228    /var/log/gdm
176    /var/log/apt
108    /var/log/ConsoleKit

```

```
64 /var/log/cups
52 /var/log/exim4
16 /var/log/samba
12 /var/log/fsck
4 /var/log/unattended-upgrades
4 /var/log/sysstat
4 /var/log/speech-dispatcher
4 /var/log/samba/cores/winbindd
4 /var/log/samba/cores
4 /var/log/news
4 /var/log/apparmor
$
```

The `-n` option allows you to sort numerically. The `-r` option lists the largest numbers first. This is perfect for finding the largest disk consumers.

The `sed` editor (see Chapters 18 and 20) will bring more clarity to this listing. To focus on the top ten disk space consumers, when line 11 is reached, `sed` is set to delete the rest of the listing. The next step is to add a line number for each line in the listing. Chapter 18 showed you how to accomplish this by adding an equal sign (=) to the `sed` command. To get those line numbers on the same line as the disk space text, combine the text lines using the `N` command, just as we did in Chapter 20. The `sed` commands needed look like this:

```
sed '{11,$D; =}' |
sed 'N; s/\n/ /' |
```

Now the output can be cleaned up using the `gawk` command (see Chapter 21). The output from the `sed` editor is piped into the `gawk` command and then printed out using the `printf` function.

```
gawk '{printf $1 ":" "\t" $2 "\t" $3 "\n"}'
```

After the line number, a colon (:) is added, and tab (`\t`)

characters are put between the individual fields for each text line's output row. This produces a nicely formatted listing of the top ten disk space consumers.

```
$ du -S /var/log |
> sort -rn |
> sed '{11,$D; =}' |
> sed 'N; s/\n/ /' |
> gawk '{printf $1 " : " "\t" $2 "\t" $3 "\n"}'
1:      6148      /var/log/dist-upgrade
2:      4864      /var/log
3:      1248      /var/log/installer
4:      1048      /var/log/dist-upgrade/20101011-1337
5:       228      /var/log/gdm
6:       176      /var/log/apt
7:       108      /var/log/ConsoleKit
8:        64      /var/log/cups
9:        52      /var/log/exim4
10:      16       /var/log/samba
$
```

Now you're in business! The next step is to use this information to create the script.

Creating the Script

To save time and effort, the script creates a report for multiple designated directories. A variable to accomplish this called `CHECK_DIRECTORIES` is used. For our purposes here, the variable is set to just two directories:

```
CHECK_DIRECTORIES="/var/log /home"
```

The script contains a `for` loop to perform the `du` command on each directory listed in the variable. This technique is used (see Chapter 12) to read and process values in a list. Each time the `for` loop iterates through the list of values in the variable `CHECK_DIRECTORIES`, it assigns to the `DIR_CHECK` variable the next value in the list:

```

for DIR_CHECK in $CHECK_DIRECTORIES
do
...
    du -S $DIR_CHECK
...
done

```

To allow quick identification, a date stamp is added to the report's file name, using the `date` command. Using the `exec` command (see Chapter 14) the script redirects its output to the date stamped report file:

```

DATE=$(date '+%m%d%y')
exec > disk_space_${DATE}.rpt

```

Now to produce a nicely formatted report, the script uses the `echo` command to put in a few report titles:

```

echo "Top Ten Disk Space Usage"
echo "for $CHECK_DIRECTORIES Directories"

```

So let's take a look at what this script looks like all put together:

```

#!/bin/bash
#
# Big_Users - find big disk space users in various directories
#####
# Parameters for Script
#
CHECK_DIRECTORIES="/var/log /home"      #directories to check
#
##### Main Script #####
#
DATE=$(date '+%m%d%y')                  #Date for report file
#
exec > disk_space_${DATE}.rpt            #Make report file Std Output
#
echo "Top Ten Disk Space Usage"          #Report header for whole report
echo "for $CHECK_DIRECTORIES Directories"
#
for DIR_CHECK in $CHECK_DIRECTORIES      #loop to du directories
do
    echo ""
    echo "The $DIR_CHECK Directory:"      #Title header for each directory
#
# Create a listing of top ten disk space users

```

```

du -S $DIR_CHECK 2>/dev/null |
sort -rn |
sed '{11,$D; =}' |
sed 'N; s/\n/ /' |
gawk '{printf $1 ":" "\t" $2 "\t" $3 "\n"}'
#
done                                     #end of for loop for du directories
#

```

And there you have it. A simple shell script that will create a date stamped report of the top ten disk space consumers for each directory you choose.

Running the Script

Before having the `Big_Users` script run automatically, you'll want to test it out a few times manually to ensure that it does what you think it should do:

```

$
$ ./Big_Users
$
$ cat disk_space_012311.rpt
Top Ten Disk Space Usage
for /var/log /home Directories

The /var/log Directory:
1:      6148      /var/log/dist-upgrade
2:      4892      /var/log
3:      1248      /var/log/installer
4:      1048      /var/log/dist-upgrade/20101011-1337
5:       176      /var/log/apt
6:       108      /var/log/ConsoleKit
7:        64      /var/log/cups
8:        52      /var/log/exim4
9:        16      /var/log/samba
10:      12       /var/log/fsck

```

```

The /home Directory:
1:    92365332    /home/user/.VirtualBox/HardDisks
2:    18659720    /home/user/Downloads
3:    17626092    /home/user/archive
4:    6174408     /home/Timothy/Junk/More_Junk
5:    6174408     /home/Timothy/Junk

```

```
6174408      /home/consultant/Work
7:           6174408      /home/consultant/Downloads
8:           3227768      /home/user/vmware/Mandriva
9:           3212464      /home/user/vmware/Fedora
10:          104632       /home/user/vmplayer
$
```

It worked! Now you can set up the shell script to execute automatically as needed. You do this using the `cron` table (see Chapter 15).

Note

When you run a script containing bash shell commands that need root privileges, you will need to use the `su` or `sudo` command. Otherwise, the script will produce unexpected results.

How often you need to run this script depends upon how active your file server is. To run the script once a week, add the following entry to your `cron` table:

```
15 7 * * 1 /home/user/Big_Users
```

This entry runs the shell script every Monday morning at 7:15 a.m. Now, you can have your coffee and review your weekly disk consumption report first thing Monday morning.

Performing Backups

Whether you're responsible for a Linux system in a business environment or just using it at home, the loss of data can be catastrophic. To help prevent bad things from happening, it's always a good idea to perform regular

backups.

However, what's a good idea and what's practical are often two separate things. Trying to arrange a backup schedule to store important files can be a challenge. This is another place where shell scripts often come to the rescue.

This section demonstrates two different methods for using shell scripts to back up data on your Linux system.

Archiving Data Files

If you're using your Linux system to work on an important project, you can create a shell script that automatically takes snapshots of specific directories. Designating these directories in a configuration file will allow you to change them when a particular project changes. This will help avoid a time-consuming restore process from your main archive files.

This section shows you how to create an automated shell script that can take snapshots of specified directories and keep an archive of past versions of your data.

The Required Functions

The workhorse for archiving data in the Linux world is the `tar` command (see Chapter 4). The `tar` command is used to archive entire directories into a single file. Here's an example of creating an archive file of a working directory using the `tar` command:

```
$ tar -cf archive.tar /home/user/backup_test
```



```
tar: Removing leading '/' from member names
$
```

The `tar` command responds with a warning message that it's removing the leading forward slash from the pathname to convert it from an absolute pathname to a relative pathname (see Chapter 3). This allows you to extract the `tar` archived files anywhere you want in your filesystem. You'll probably want to get rid of that message in your script. Do that by redirecting `STDERR` to the `/dev/null` file (see Chapter 14):

```
$ tar -cf archive.tar /home/user/backup_test 2>/dev/null
$
```

Because a `tar` archive file can consume a lot of disk space, it's a good idea to compress the file. You can do this by simply adding the `-z` option. This will compress the `tar` archive file into a gzipped `tar` file, which is called a tarball. Be sure to use the proper file extensions to denote that the file is a tarball. Either `.tar.gz` or `.tgz` is fine. Here's an example of creating a tarball of a working directory:

```
$ tar -zcf archive.tar.gz /home/user/backup_test 2>/dev/null
$
```

Now you have the main component for your archive script completed.

Instead of modifying or creating a new archive script for each new directory or file you want to back up, you can use a configuration file. The configuration file should contain each directory or file you want to be included in the archive:

```
$ cat Files_To_Backup
/home/user/Downloads
/home/user/Documents/CLandSS_V2
/home/Samantha/Documents
```

```
/home/Does_not_exist
/home/Timothy/Junk
/home/consultant/Work
$
```

Caution

If you're using a Linux distribution that includes a graphical desktop, be careful about archiving your entire `$HOME` directories. While this may be tempting, the `$HOME` directory contains lots of configuration and temporary files related to the graphical desktop. It will create a much larger archive file than you probably intended. Pick a subdirectory in which to store your working files, and use that subdirectory in your archive configuration file.

You can have the script read through the configuration file and add the names of each directory or file to the archive list. To do this, use the simple `read` command (see Chapter 13) to read each record from the file. But instead of using the `cat` command piped into a `while` loop, as we did in Chapter 13, this script is going to redirect `STDINPUT` using the `exec` command (see Chapter 14). Here's how it will look:

```
exec < $CONFIG_FILE
```

```
read FILE_NAME
```

Notice that a variable is used for the archive configuration file as well as each record read in from the file. As long as the `read` command finds a new record to read in the configuration file, it will return an exit value of 0 for success in the `?` variable (see Chapter 10). You can use this as a test in a `while` loop in order to read all the records from the configuration file:

```
while [ $? -eq 0 ]
do
...
read FILE_NAME
done
```

Once the `read` command hits the end of the configuration file, it returns a non-successful status. At that point, the `while` loop is exited.

In the `while` loop, two things need to happen. First, you must add the file or directory name to your archive list. Yet even more important is to check and see if that file or directory even exists! It would be very easy to remove a directory from the filesystem and forget to update the archive configuration file. You can check a file's or directory's existence using a simple `if` statement and the test command file comparisons (see Chapter 11). If the file or directory does exist, it is added to the list of files to archive, `FILE_LIST`. Otherwise, a warning message is issued. Here is what this looks like:

```
if [ -f $FILE_NAME -o -d $FILE_NAME ]
then
    # If file exists, add it's name to the list.
    FILE_LIST="$FILE_LIST $FILE_NAME"
else
    # If file doesn't exists, issue warning
    echo
    echo "$FILE_NAME, does not exist."
    echo "Obviously, I will not include it in this archive."
    echo "It is listed on line $FILE_NO of the config file."
    echo "Continuing to build archive list..."
    echo

fi

#
FILE_NO=$((FILE_NO + 1)) # Increase Line/File number by one.
```

Because a record in our archive configuration file can be a file name or a directory, the `if` statement tests for the

existence of both, using the `-f` and the `-d` options. The `or` option, `-o`, allows for either the file's or the directory's existence test to return a `true` for the entire `if` statement to be treated as `true`.

To provide a little extra help in tracking down non-existent directories and files, the variable `FILE_NO` is added. Thus, the script can tell you exactly what line number in the archive configuration file contains the incorrect or missing file or directory.

You should now have enough information to start building the script. The next section walks you through creating the archive script.

Creating a Daily Archive Script

The `Daily_Archive` script automatically creates an archive to a designated location, using the current date to uniquely identify the file. Here's the code for that portion of the script:

```
DATE='date +%y%m%d'
#
# Set Archive File Name
#
FILE=archive$DATE.tar.gz
#
# Set Configuration and Destination File
#
CONFIG_FILE=/home/user/archive/Files_To_Backup
DESTINATION=/home/user/archive/$FILE
#
```

The `DESTINATION` variable appends the full pathname for the archived file. The `CONFIG_FILE` variable points to the archive configuration file containing the files to be archived. These both can be easily changed to alternate directories and

files if needed.

The `Daily_Archive` script, all put together, now looks like this:

```
#!/bin/bash
#
# Daily_Archive - Archive designated files & directories
#####
#
# Gather Current Date
#
DATE=`date +%y%m%d`
#
# Set Archive File Name
#
FILE=archive$DATE.tar.gz
#
# Set Configuration and Destination File
#
CONFIG_FILE=/home/user/archive/Files_To_Backup
DESTINATION=/home/user/archive/$FILE
#
##### Main Script #####
#
# Check Backup Config file exists
#
if [ -f $CONFIG_FILE ]      # Make sure the config file still exists.
then                        # If it exists, do nothing but continue on.
    echo
else                        # If it doesn't exist, issue error & exit script.
    echo
    echo "$CONFIG_FILE does not exist."
    echo "Backup not completed due to missing Configuration File"
    echo
    exit
fi
#
# Build the names of all the files to backup
#
FILE_NO=1                  # Start on Line 1 of Config File.
exec < $CONFIG_FILE        # Redirect Std Input to name of Config File
#
read FILE_NAME              # Read 1st record
#
while [ $? -eq 0 ]          # Create list of files to backup.
do
    # Make sure the file or directory exists.
    if [ -f $FILE_NAME -o -d $FILE_NAME ]
```

```

then
    # If file exists, add its name to the list.
    FILE_LIST="$FILE_LIST $FILE_NAME"
else
    # If file doesn't exist, issue warning
    echo
    echo "$FILE_NAME, does not exist."
    echo "Obviously, I will not include it in this archive."
    echo "It is listed on line $FILE_NO of the config file."
    echo "Continuing to build archive list..."
    echo
fi

#
FILE_NO=$((FILE_NO + 1)) # Increase Line/File number by one.
read FILE_NAME           # Read next record.

done
#
#####
#
# Backup the files and Compress Archive
#
tar -czf $DESTINATION $FILE_LIST 2> /dev/null
#

```

Running the Daily Archive Script

Testing the `Daily_Archive` script is straightforward:

```

$ ./Daily_Archive

/home/Does_not_exist, does not exist.
Obviously, I will not include it in this archive.
It is listed on line 4 of the config file.
Continuing to build archive list...

```

You can see that the script caught one directory that does not exist, `/home/Does_not_exist`. It lets you know what line number in the configuration file this erroneous directory is on, and then continues on making a list and archiving the data. The data is now safely archived in a tarball file.

Creating an Hourly Archive Script

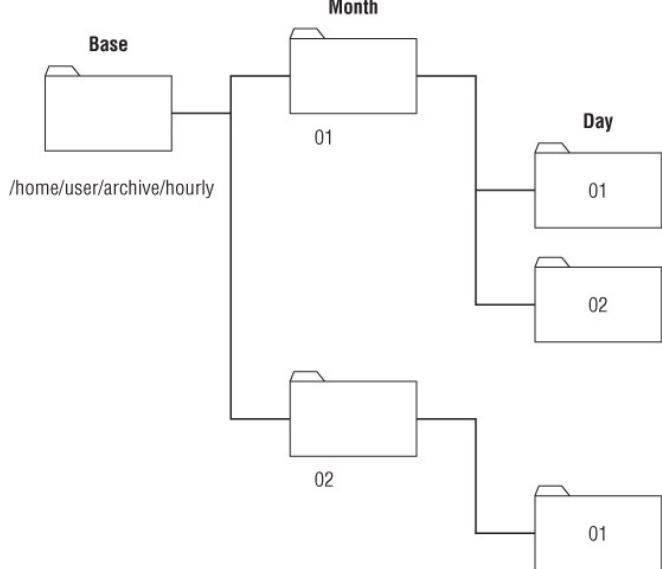
If you are in a high-volume production environment where files are changing rapidly, a daily archive might not be good enough. If you want to increase the archiving frequency to hourly, you'll need to take another item into consideration.

When backing up files hourly and trying to use the date command to timestamp each file name, things can get pretty ugly pretty quickly. Sifting through a directory of files with file names looking like this is tedious:

```
archive010211110233.tar.gz
```

Instead of placing all of the archive files in the same folder, you'd be better off creating a directory hierarchy for your archived files. [Figure 26.1](#) demonstrates this principle.

[Figure 26.1](#) Creating an archive directory hierarchy



The archive directory contains directories for each month of the year, using the month number as the directory name. Each month's directory in turn contains folders for each day of the month (using the day's numerical value as the directory name). This allows you to just timestamp the individual archive files and then place them in the appropriate directory for the day and month.

Now there is a new challenge to solve. The script must create the individual month and day directories

automatically and know that if they already exist, they don't need to be created.

If you peruse the command line options for the `mkdir` command (see Chapter 3), you'll find the `-p` command line option. This option allows you to create directories and subdirectories in a single command; plus, the added benefit is that it doesn't produce an error message if the directory already exists. Perfect!

We're now ready to create the `Hourly_Archive` script. Here is the top half of the script:

```
#!/bin/bash
#
# Hourly_Archive - Every hour create an archive
#####
#
# Set Configuration File
#
CONFIG_FILE=/home/user/archive/hourly/Files_To_Backup
#
# Set Base Archive Destination Location
#
BASEDEST=/home/user/archive/hourly
#
# Gather Current Day, Month & Time
#
DAY=`date +%d`
MONTH=`date +%m`
TIME=`date +%k%M`
#
# Create Archive Destination Directory
#
mkdir -p $BASEDEST/$MONTH/$DAY
#
# Build Archive Destination File Name
#
DESTINATION=$BASEDEST/$MONTH/$DAY/archive$TIME.tar.gz
#
##### Main Script #####
...
```

Once the script reaches the “Main Script” portion of

Hourly_Archive, the script is an exact duplicate of the Daily_Archive script. A lot of the work has already been done!

Hourly_Archive retrieves the day and month values from the date command, along with the timestamp used to uniquely identify the archive file. It then uses that information to create the archive directory for the day (or just silently exit if it already exists). Finally, the script uses the tar command to create the archive and compress it into a tarball.

Running the Hourly Archive Script

As with the Daily_Archive script, it's a good idea to test the Hourly_Archive script before putting it in the cron table:

```
$ cat /home/user/archive/hourly/Files_To_Backup
/home/Development/Simulation_Logs
/home/Production/Machine_Errors
$
$ ./Hourly_Archive
$
$ ls /home/user/archive/hourly/01/02
archive1601.tar.gz
$
$ ./Hourly_Archive
$
$ ls /home/user/archive/hourly/01/02
archive1601.tar.gz      archive1602.tar.gz
$
```

The script worked fine the first time, creating the appropriate month and day directories, and then creating the archive file. Just to test things out, it was run a second time to see if it would have a problem with the existing directories. The script again ran fine and created the second archive file. It's now ready for the cron table.

Managing User Accounts

Managing user accounts is much more than just adding, modifying, and deleting accounts. You also have to consider security issues, the need to preserve work, and the accurate management of the accounts. This can be a time-consuming task. Here is another instance when writing script utilities is a real time saver!

The Required Functions

Deleting an account is the more complicated task of managing accounts. When deleting an account, at least four separate actions are required:

1. Obtain the correct user account name to delete.
2. Kill any processes currently running on the system that belongs to that account.
3. Determine all files on the system belonging to the account.
4. Remove the user account.

It's easy to miss a step. The shell script utility in this section will help you avoid making such mistakes.

Obtain Correct Account Name

The first step in the account deletion process is the most important, obtaining the correct user account name to delete. Because this is an interactive script, you can use

the `read` command (see Chapter 13) to obtain the account name. If the script user walks away and leaves the question hanging, you can use the `-t` option on the `read` command and timeout after giving the script user 60 seconds to answer the question:

```
echo "Please enter the username of the user "  
echo -e "account you wish to delete from system: \c"  
read -t 60 ANSWER
```

Because interruptions are part of life, it's best to give users three chances to answer the question. To accomplish this, use a `while` loop (see Chapter 12) with the `-z` option to test if the `ANSWER` variable is empty. The `ANSWER` variable will be zero-length when the script first enters the `while` loop because the question will be put at the end of the loop:

```
while [ -z "$ANSWER" ]  
do  
...  
echo "Please enter the username of the user "  
echo -e "account you wish to delete from system: \c"  
read -t 60 ANSWER  
done
```

Now you need a way to communicate with the script user about when the first question timeout occurred, when there is one more chance to answer the question, and so on. The `case` statement (see Chapter 11) is the structured command that will work perfectly here. Using the incremented `ASK_COUNT` variable, you can set up different messages to communicate to the script user. The code for this section looks like this:

```
case $ASK_COUNT in  
2)  
    echo  
    echo "Please answer the question."  
    echo
```

```

;;
3)

    echo
    echo "One last try...please answer the question."
    echo

;;
4)

    echo
    echo "Since you refuse to answer the question..."
    echo "exiting program."
    echo
    #
    exit

;;
esac
#

```

Now the script has all the structure it needs to ask the user what account to delete. There are several more questions in this script you will need to ask the user and just asking that one question was a lot of code! Therefore, let's turn this piece of code into a function (see Chapter 16) in order to use it in multiple locations in your `Delete_User` script.

The first thing you need to do is to declare the function's name, `get_answer`. Next, clear out any previous answers to questions your script user gave using the `unset` command (see Chapter 5). The code to do these two items looks like this:

```

function get_answer {
#
unset ANSWER

```

The other item you need to change in the original code is the actual question to the script user. The script won't be asking the same question each time, so let's create two new variables, `LINE1` and `LINE2`, to handle question lines:

```

echo $LINE1
echo -e $LINE2" \c"

```

However, not every question has two lines to display. Some will have only one line. Use an `if` construct (see Chapter 11) to assist with this problem. The function will test if `LINE2` is empty and only use `LINE1` if it is:

```
if [ -n "$LINE2" ]
then
    echo $LINE1
    echo -e $LINE2" \c"
else
    echo -e $LINE1" \c"
fi
```

Finally, our function will need to clean up after itself by clearing out the `LINE1` and `LINE2` variables. Thus, the function now looks like this:

```
function get_answer {
#
unset ANSWER
ASK_COUNT=0
#
while [ -z "$ANSWER" ]
do
    ASK_COUNT=$(( ASK_COUNT + 1 ))
#
    case $ASK_COUNT in
        2)
            echo
...
    esac
#
    echo
    if [ -n "$LINE2" ]
    then
        #Print 2 lines
        echo $LINE1
        echo -e $LINE2" \c"
    else
        #Print 1 line
        echo -e $LINE1" \c"
    fi
#
    read -t 60 ANSWER
done
#
unset LINE1
unset LINE2
```

```
} #End of get_answer function
```

To ask the script user what account to delete, you need to set a few variables and then call the `get_answer` function. Using the new function makes the script code much simpler:

```
LINE1="Please enter the username of the user "  
LINE2="account you wish to delete from system:"  
get_answer  
USER_ACCOUNT=$ANSWER
```

Because of potential typographical errors, you will want to verify the user account name that was entered. This will be easy because we already have the code in place to handle asking a question:

```
LINE1="Is $USER_ACCOUNT the user account "  
LINE2="you wish to delete from the system? [y/n]"  
get_answer
```

Once the question is asked, the script needs to process the answer. The variable `ANSWER` is once again carrying the script user's answer to the question. If the user answered "yes," you have the correct user account to delete and can continue on in the script. You can use the `case` statement (see Chapter 11) to process the answer. Be sure to code the `case` statement so it will check for the multiple ways the answer "yes" can be entered.

```
case $ANSWER in  
y|Y|YES|yes|Yes|yEs|yeS|YEs|yES )  
#  
;;  
*)  
  
    echo  
    echo "Because the account, $USER_ACCOUNT, is not "  
    echo "the one you wish to delete, we are leaving the script..."  
    echo  
    exit
```

```
;;  
esac
```

There will be many times the script will need to handle a yes/no answer from the user. Thus, once again, it makes sense to create a function to handle this task. Only a few changes need to be made to the preceding code. Simply declare the function's name and add the variables `EXIT_LINE1` and `EXIT_LINE2` to the `case` statement. These changes, along with some variable cleanup at the end, will result in the `process_answer` function:

```
function process_answer {  
#  
case $ANSWER in  
y|Y|YES|yes|Yes|yEs|yeS|YEs|yES )  
;;  
*)  
    echo  
    echo $EXIT_LINE1  
    echo $EXIT_LINE2  
    echo  
    exit  
;;  
esac  
#  
unset EXIT_LINE1  
unset EXIT_LINE2  
#  
} #End of process_answer function
```

A simple function call will now process the answer:

```
EXIT_LINE1="Because the account, $USER_ACCOUNT, is not "  
EXIT_LINE2="the one you wish to delete, we are leaving the script..."  
process_answer
```

The user has given us the name of the account to delete and has verified it. Now is a good time to double check that the user account really exists on the system. Also, it is a good idea to show the full account record to the script user to double check this is really the account to delete. To

accomplish these items, use a variable, `USER_ACCOUNT_RECORD`, and set it to the outcome of a `grep` (see Chapter 4) search through the `/etc/passwd` file for the user account. The `-w` option will allow you to make an exact word match in the user account records of this particular user account:

```
USER_ACCOUNT_RECORD=$(cat /etc/passwd | grep -w $USER_ACCOUNT)
```

If no user account record is found in `/etc/passwd`, that means the account has already been deleted or never existed in the first place. In either case, you need to let the script user know this and exit the script. The exit status of the `grep` command will help us out here. If the account record is not found, the `?` variable will be set to 1:

```
if [ $? -eq 1 ]
then
    echo
    echo "Account, $USER_ACCOUNT, not found. "
    echo "Leaving the script..."
    echo
    exit
fi
```

If the record was found, you still need to verify with the script user that this is the correct account. Here is where all the work we did to set up our functions will really pay off! All you will have to do is set the proper variables and call the functions:

```
echo "I found this record:"
echo $USER_ACCOUNT_RECORD
echo
#
LINE1="Is this the correct User Account? [y/n]"
get_answer
#
EXIT_LINE1="Because the account, $USER_ACCOUNT, is not"
EXIT_LINE2="the one you wish to delete, we are leaving the script..."
process_answer
```

Remove Account Processes

So far, you have obtained and verified the correct name of the user account to be deleted. In order to remove the user account from the system, the account cannot own any processes currently running. Thus, your next step is to find and kill off those processes. This is going to get a little complicated!

Finding the user processes will be the easy part. Here the script can use the `ps` command (see Chapter 4) and the `-u` option to locate any running processes owned by the account. If any processes do exist, they will be displayed to the script user.

```
ps -u $USER_ACCOUNT
```

You can use the `ps` command's exit status and a `case` structure to determine the next step to take:

```
case $? in
1)
    echo "There are no processes for this account currently running."
    echo
;;
0)
    unset ANSWER
    LINE1="Would you like me to kill the process(es)? [y/n]"
    get_answer
    ...
esac
```

If the exit status returns a 1, then there are no processes running on the system that belong to the user account. However, if the exit status returns a 0, then there are processes owned by this account running on the system. In this case, the script needs to ask our script user if they

would like to have these processes killed. You can accomplish this task by using the `get_answer` function.

You might think that the next action the script will do is to call the `process_answer` function. Unfortunately, the next item is too complicated for `process_answer`. You need to embed another `case` statement to process the script user's answer. The first part of the `case` statement will look very similar to our `process_answer` function:

```
case $ANSWER in
    y|Y|YES|yes|Yes|yEs|yeS|YES|yES ) # If user answers "yes",
                                          #kill User Account processes.
        ...
    ;;
    *) # If user answers anything but "yes", do not kill.
        echo
        echo "Will not kill the process(es)"
        echo
    ;;
esac
```

As you can see, there is nothing too interesting in the `case` statement itself. Where things do get interesting is within the "yes" section of the `case` statement. Here, the user account processes need to be killed. To do so, reissue the `ps` command, this time sending the output to a temporary report file instead of the screen.

```
ps -u $USER_ACCOUNT > $USER_ACCOUNT_Running_Process.rpt
```

You can use the `exec` command and a `while` loop to read through the report file, similar to the method we used to read through the backup configuration file in the `Daily_Archive` script.

```
exec < $USER_ACCOUNT_Running_Process.rpt
read USER_PROCESS_REC
while [ $? -eq 0 ]
do
```

```
...
read USER_PROCESS_REC
done
```

Within the `while` loop, extract the process ID (PID) from each running process's status record using the `cut` command. Once you have the PID, use the `kill` command (see Chapter 4) with a `-9` option to unconditionally terminate the process.

```
USER_PID=$(echo $USER_PROCESS_REC | cut -d " " -f1)
kill -9 $USER_PID
echo "Killed process $USER_PID"
```

Now with any user account–owned processes killed, the script can move on to the next step, finding all of the user account's files.

Find Account Files

When a user account is deleted from the system, it is a good practice to archive all the files that belonged to that account. Along with that practice, it is also important to remove the files or assign their ownership to another account. If the account you delete has a User ID of 1003, and you don't remove or reassign those files, then the next account that is created with a User ID of 1003 will own those files! You can see the security disasters that can occur in this scenario.

The `Delete_User` script will not do all of that for you, but it will create a report that can be used in the `Daily_Archive` script as a backup configuration file. And you can use the report to help you remove or reassign the files.

To find the user's files, you can use the `find` command. The `find` command searches the entire filesystem with the `-u` option, which will pinpoint any user account–owned files. The command looks like the following:

```
find / -user $USER_ACCOUNT > $REPORT_FILE
```

That was pretty simple compared to dealing with the user account processes! It gets even easier in the next step of the `Delete_User` script, actually removing the user account.

Remove Account

It's always a good idea to be a little paranoid about removing a user account from the system. Therefore, you should ask one more time if the script user really wants to remove the account.

```
LINE1="Do you wish to remove $User_Account's account from system? [y/n]"
get_answer
#
EXIT_LINE1="Since you do not wish to remove the user account,"
EXIT_LINE2="$USER_ACCOUNT at this time, exiting the script..."
process_answer
```

Finally, we get to the main purpose of our script, actually removing the user account from the system. Here you can use the `userdel` command (see Chapter 6):

```
userdel $USER_ACCOUNT
```

Now that we've got all the pieces, we are ready to put them together into a whole, useful script utility.

Creating the Script

Remember that the `Delete_User` script is highly interactive

with the script's user. Therefore, it is important to include a lot of verbiage to keep the user informed about what is going on during the script's execution.

At the top of the script, the two functions `get_answer` and `process_answer` are declared. The script then steps into the four steps of removing the user: obtaining and confirming the user account name; finding and killing the user's processes; creating a report of all files owned by the user account; and actually removing the user account.

Here's the entire `Delete_User` script:

```
#!/bin/bash
#
#Delete_User - Automates the 4 steps to remove an account
#
#####
# Define Functions
#
#####
function get_answer {
#
unset ANSWER
ASK_COUNT=0
#
while [ -z "$ANSWER" ]      #While no answer is given, keep asking.
do
    ASK_COUNT=$(( ASK_COUNT + 1 ])
#
    case $ASK_COUNT in      #If user gives no answer in time allotted
        2)
            echo
            echo "Please answer the question."
            echo
            ;;
        3)
            echo
            echo "One last try...please answer the question."
            echo
            ;;
        4)
            echo
            echo "Since you refuse to answer the question..."
            echo "exiting program."
```

```

        echo
        #
        exit
    ;;
esac

#
echo

#

if [ -n "$LINE2" ]
then
    #Print 2 lines
    echo $LINE1
    echo -e $LINE2" \c"
else
    #Print 1 line
    echo -e $LINE1" \c"
fi

#

#    Allow 60 seconds to answer before time-out
read -t 60 ANSWER

done

# Do a little variable clean-up
unset LINE1
unset LINE2
#
} #End of get_answer function
#
#####
function process_answer {
#
case $ANSWER in
y|Y|YES|yes|Yes|yEs|yeS|YEs|yES )
# If user answers "yes", do nothing.
;;
*)
# If user answers anything but "yes", exit script
echo
echo $EXIT_LINE1
echo $EXIT_LINE2
echo
exit

;;
esac

#
# Do a little variable clean-up
#
unset EXIT_LINE1
unset EXIT_LINE2
#
} #End of process_answer function
#
#####

```

End of Function Definitions

```
#
##### Main Script #####
# Get name of User Account to check
#
echo "Step #1 - Determine User Account name to Delete "
echo
LINE1="Please enter the username of the user "
LINE2="account you wish to delete from system:"
get_answer
USER_ACCOUNT=$ANSWER
#
# Double check with script user that this is the correct User Account
#
LINE1="Is $USER_ACCOUNT the user account "
LINE2="you wish to delete from the system? [y/n]"
get_answer
#
# Call process_answer function:
#     if user answers anything but "yes", exit script
#
EXIT_LINE1="Because the account, $USER_ACCOUNT, is not "
EXIT_LINE2="the one you wish to delete, we are leaving the script..."
process_answer
#
#####
# Check that USER_ACCOUNT is really an account on the system
#
USER_ACCOUNT_RECORD=$(cat /etc/passwd | grep -w $USER_ACCOUNT)
#
if [ $? -eq 1 ]           # If the account is not found, exit script
then
    echo
    echo "Account, $USER_ACCOUNT, not found. "
    echo "Leaving the script..."
    echo
    exit
fi
#
echo
echo "I found this record:"
echo $USER_ACCOUNT_RECORD
echo
#
LINE1="Is this the correct User Account? [y/n]"
get_answer
#
#
# Call process_answer function:
#     if user answers anything but "yes", exit script
```



```

EXIT_LINE1="Because the account, $USER_ACCOUNT, is not "
EXIT_LINE2="the one you wish to delete, we are leaving the script..."
process_answer
#
#####
# Search for any running processes that belong to the User Account
#
echo
echo "Step #2 - Find process on system belonging to user account"
echo
echo "$USER_ACCOUNT has the following processes running: "
echo
#
ps -u $USER_ACCOUNT          #List user processes running.

case $? in
1)      # No processes running for this User Account
        #
        echo "There are no processes for this account currently running."
        echo
;;
0)      # Processes running for this User Account.
        # Ask Script User if wants us to kill the processes.
        #
        unset ANSWER
        LINE1="Would you like me to kill the process(es)? [y/n]"
        get_answer
        #
        case $ANSWER in
y|Y|YES|yes|Yes|yEs|yeS|YES|yES )      # If user answers "yes",
                #kill User Account processes.
                #
                echo
                #
                # Clean-up temp file upon signals
                trap "rm $USER_ACCOUNT_Running_Process.rpt" SIGTERM SIGINT SIGQUIT
                #
                # List user processes running
                ps -u $USER_ACCOUNT > $USER_ACCOUNT_Running_Process.rpt
                #
                exec < $USER_ACCOUNT_Running_Process.rpt      # Make report Std Input
                #
                read USER_PROCESS_REC      # First record will be blank
                read USER_PROCESS_REC
                #
                while [ $? -eq 0 ]
                do
                        # obtain PID
                        USER_PID=$(echo $USER_PROCESS_REC | cut -d " " -f1)

```

```

kill -9 $USER_PID
echo "Killed process $USER_PID"
read USER_PROCESS_REC
done
#
echo
rm $USER_ACCOUNT_Running_Process.rpt          # Remove temp report.
;;
*)      # If user answers anything but "yes", do not kill.
echo
echo "Will not kill the process(es)"
echo

;;
esac

;;
esac

;;
esac

#####
# Create a report of all files owned by User Account
#
echo
echo "Step #3 - Find files on system belonging to user account"
echo
echo "Creating a report of all files owned by $USER_ACCOUNT."
echo
echo "It is recommended that you backup/archive these files,"
echo "and then do one of two things:"
echo "  1) Delete the files"
echo "  2) Change the files' ownership to a current user account."
echo
echo "Please wait. This may take a while..."
#
REPORT_DATE='date +%y%m%d'
REPORT_FILE=$USER_ACCOUNT"_Files_"$REPORT_DATE".rpt"
#
find / -user $USER_ACCOUNT > $REPORT_FILE 2>/dev/null
#
echo
echo "Report is complete."
echo "Name of report:      $REPORT_FILE"
echo "Location of report:   'pwd'"
echo
#####
# Remove User Account
echo
echo "Step #4 - Remove user account"
echo
#
LINE1="Do you wish to remove $User_Account's account from system? [y/n]"
get_answer

```

```
# Call process_answer function:
#         if user answers anything but "yes", exit script
#
EXIT_LINE1="Since you do not wish to remove the user account,"
EXIT_LINE2="$USER_ACCOUNT at this time, exiting the script..."
process_answer
#
userdel $USER_ACCOUNT          #delete user account
echo
echo "User account, $USER_ACCOUNT, has been removed"
echo
#
```

That was a lot of work! However, the `Delete_User` script will be a great timesaver and will help you avoid a lot of nasty problems when deleting user accounts.

Running the Script

Because it is intended to be an interactive script, the `Delete_User` script should not be placed in the `cron` table. However, it is still important to ensure it works as expected. We'll test the script by removing a consultant account that was set up for a temporary consultant on the system.

```
$ ./Delete_User
Step #1 - Determine User Account name to Delete

Please enter the username of the user
account you wish to delete from system: consultant

Is consultant the user account
you wish to delete from the system? [y/n] y

I found this record:
consultant:x:1004:1001:Consultant,,,:/home/consultant:/bin/bash

Is this the correct User Account? [y/n] y

Step #2 - Find process on system belonging to user account

consultant has the following processes running:
```

```
PID TTY TIME CMD
There are no processes for this account currently running.

Step #3 - Find files on system belonging to user account

Creating a report of all files owned by consultant.

It is recommended that you backup/archive these files,
and then do one of two things:
1) Delete the files
2) Change the files' ownership to a current user account.

Please wait. This may take a while...

Report is complete.
Name of report: consultant_Files_110123.rpt
Location of report: /home/Christine

Step #4 - Remove user account

Do you wish to remove consultant's account from system? [y/n] y

User account, consultant, has been removed
$
```

That worked great! Now you have a script utility that will be of assistance when you need to delete user accounts. Even better is that you can modify it to meet your organization's needs.

Summary

This chapter put some of the shell-scripting information presented in the book to good use for creating Linux utilities. When you're responsible for a Linux system, whether it's a large multi-user system or your own system, there are lots of things you need to watch. Instead of manually running commands, you can create shell script

utilities to do the work for you.

The chapter demonstrated how to use the `du` command to determine disk space consumption. The `sed` and `gawk` commands were then used to retrieve specific information from the data. Passing the output from a command to `sed` and `gawk` to parse data is a common function in shell scripts, so it's a good idea to know how to do it.

The next section walked you through using shell scripts for archiving and backing up data files on the Linux system. The `tar` command is a popular command for archiving data. The chapter showed you how to use it in shell scripts to create archive files, and how to manage the archive files in an archive directory.

We ended the chapter with information on using a shell script for the four steps needed to delete user accounts. Creating functions for shell code that is repeated within a script makes the code easier to read and modify. This script combined many of the different structured commands, such as the `case` and `while` commands. The chapter demonstrated the difference in script structure for a script destined for the `cron` tables versus an interactive script.

Now we will move on to more advanced shell scripts that will help with more complex Linux system issues.

Chapter 27

Advanced Shell Scripts

In This Chapter

- Monitoring system statistics
- Problem-tracking database

In this chapter, we will look at more advanced methods to use in shell scripts. These techniques enable you to work with scripts on your system in a variety of different ways.

Monitoring System Statistics

There are many utilities in the bash shell to help you keep an eye on Linux system performance and resource usage. Unfortunately, it is difficult to find the time to use them, especially when you have several systems to manage. Creating a few shell scripts to create reports using the output from these utilities can be an excellent use of your time.

This section will take you through the creation of some advanced shell scripts to assist you in monitoring various

system performance statistics.

System Snapshot Report

A nice way to stay informed about your system's performance and resource usage is via a snapshot report. A *snapshot report* is a picture of a system's statistics at a certain point in time. A report of this kind is an “executive summary” of a system's health.

The script to create this snapshot report can be run as many times a day as you want a picture. Even better, the script will send you the report via e-mail.

The Required Functions

You will need to use four different bash shell commands to create the snapshot report: `uptime`, `df`, `free`, and `ps`. These various commands will give you the statistics needed.

Up Time

The first command, `uptime`, is the most basic system statistics command:

```
$ uptime
14:15:23 up 1 day, 5:10, 3 users, load average: 0.66, 0.54, 0.33
```

The `uptime` command gives a few different basic pieces of information that we can use:

- The current time
- The number of days, hours, and minutes the system

has been operational

- The number of users currently logged into the system
- The one, five, and fifteen minute load averages

The system's operational time is a great statistic to grab for our snapshot report. However, there is a problem in trying to grab this statistic. Your system may have only been up for a few minutes or it could have been up for days. (Rumor has it that there are even some Linux servers that have run for years!) Here's the `uptime` command output for a system that has been up for a few hours:

```
$ uptime
18:00:20 up 8:55, 3 users, load average: 0.62, 0.45, 0.37
```

And here's the `uptime` output for a system that has been up for a few days:

```
$ uptime
13:29:32 up 3 days, 4:24, 4 users, load average: 1.44, 0.83, 0.46
```

As you can see, it's going to be tricky to gather the right statistic from this listing. The `gawk` command (see Chapter 16) can do this job with no problems for a system that has been up for a few days:

```
$ uptime | gawk '{print $2,$3,$4,$5}'
up 3 days, 4:27,
```

However, this is what happens when the same `gawk` command is used on a system that has only been up for a few hours:

```
$ uptime | gawk '{print $2,$3,$4}'
up 8:59, 3 users,
```

To solve the dilemma, you can use a more advanced feature of `gawk`, the `if` statement (see Chapter 19). If the

system has been up for several days, the `gawk` variable `$4` will contain the word “days.” To test `$4` for the string “days,” your `gawk` command should start to look like this:

```
$ uptime |
> gawk '{if ($4 == "days") {print $2,$3,$4,$5} else {print $2,$3}}'
up 9:23,
```

You may run into the situation where the system has only been up for a single day. Thus, you need to test for the string “day” as well. The Boolean operator `or` (`||`) will do the trick. Thus, your `gawk` command now looks like this:

```
gawk '{if ($4 == "days" || $4 == "day") {print $2,$3,$4,$5} else {print $2,$3}}'
```

Now what about those unsightly commas in the output? Using the `sed` command, you can simply replace them with spaces.

```
$ uptime | sed -n '/,/s/,/ /gp' |
> gawk '{if ($4 == "days" || $4 == "day") {print $2,$3,$4,$5}
> else {print $2,$3}}'
up 9:32
```

That looks much better! And just to be sure, let's test out the command on a system that has been up for more than one day:

```
$ uptime | sed -n '/,/s/,/ /gp' |
> gawk '{if ($4 == "days" || $4 == "day") {print $2,$3,$4,$5}
> else {print $2,$3}}'
up 3 days 4:34
```

That looks good, too. Now we can move on to another performance utility command to use for creating the system snapshot report.

Disk Usage

The next command to help you create your snapshot report

is `df`. The `df` command shows disk space usage statistics:

```
$ df
Filesystem      1K-blocks      Used      Available  Use%    Mounted on
/dev/sda2       307465076    185930336    105916348    64%     /
none            503116       240          502876       1%      /dev
none            508716       252          508464       1%      /dev/shm
none            508716       100          508616       1%      /var/run
...
```

The `df` command shows the current disk space statistics for all of the real and virtual disks on the system. You can pass an individual disk name to the `df` command to retrieve information only for that particular disk. To help make the information easier to read, add the `-h` (human readable) option to the command:

```
$ df -h /dev/sda2
Filesystem Size Used Avail Use% Mounted on
/dev/sda2  294G 178G 102G   64%   /
$
```

The `sed` and `gawk` commands will help you again in parsing and printing out the needed information for the report. Because we don't need to include the header line from the `df` command's output, use `sed` to search for something that is unique to the data lines. The string `% /` is only found in the data lines. However, `sed` recognizes a forward slash as part of its command structure. To solve this problem, put the escape character `\` before the forward slash. This will allow `sed` to see the forward slash as part of the string and conduct the search properly. Finally, use `gawk` to print the disk usage percentage, which is in the `gawk` variable `$5`. The entire command line will look like this:

```
$ df -h /dev/sda2 | sed -n '% /\|p' | gawk '{print $5}'
64%
```

This works great. Now let's move on to getting a snapshot of your system's memory.

Memory Usage

There are several different bash shell commands you can use to obtain memory statistics. For our purposes here, we are going to concentrate on the `free` command.

The `free` command displays the total amount of physical memory, along with how much of it is free and how much of it is allocated. It also displays that same data for the swap memory, and includes the kernel buffers. Here's an example of output from the `free` command:

```
$ free
      total          used          free      shared    buffers     cached
Mem:    1017436      890148      127288         0       31536      484100
-/+ buffers/cache: 374512      642924
Swap:    200776         0         200776
$
```

Because we pulled a percentage statistic for disk usage, it is a good idea to do the same for memory usage. However, you will notice that the `free` command does not provide any usage percentages! Not a problem. Using the `gawk` command, you can calculate a percentage.

The `free` command output is already nicely divided by spaces, so `gawk` can simply reference the two fields needed for the calculation by their variable names. The `$2` variable is the total memory available and the `$3` variable is the total used memory. By dividing `$3` by `$2`, `gawk` will provide you a percentage of memory used:

```
$ free | sed -n '2p' | gawk 'x = ($3 / $2) {print x}'
```

Notice that `sed` was used in the command to pull only the second line of data from the `free` command output, which allows you to “throw out” the header line. The number produced is still not quite in the format needed for the report. You can clean up the percentage via multiplying the number by 100 and truncating it with the `gawk` integer function, `int`. Finally, tack on a percentage sign using `sed`:

```
$ free | sed -n '2p' |
> gawk 'x = int(($3 / $2) * 100) (print x)' |
> sed 's/$/%/'
87%
```

Now the memory statistic is in the format you need for the report. Let's move on to the last item to include—zombies!

Zombie Processes

Zombies on a Linux system are processes caught in limbo. The process has completed its work, but is stuck in an uncompleted state for various reasons. Thus, like its namesake, a zombie process is neither dead nor alive.

While a zombie process or two is not a big deal on a Linux system, several of them can cause problems. A zombie holds a process ID, which will not be released back to the process ID table until the zombie is killed off. Thus, another system statistic to watch is the current number of zombie processes.

This is easily done with the `ps` command (see Chapter 4). The `ps` command along with the `-al` options will return all the processes on the system along with their process states. A

state of `Z` indicates the process is a zombie:

```
$ ps -al
F S  UID      PID PPID C   PRI NI ADDR      SZ WCHAN  TTY          TIME CMD
0 T 1000    2174 2031 0   80  0   - 2840 signal pts/0 00:00:00 mail
0 T 1000    2175 2031 0   80  0   - 2839 signal pts/0 00:00:00 mail
1 Z 1000    8779 8797 0   80  0   -    0 exit   pts/0 00:00:00 zomb <defunct>
0 R 1000    8781 8051 0   80  0   - 1094 -      pts/5 00:00:00 ps
$
```

Once again, the `gawk` command will help you here. With the output from the `ps` command piped into `gawk`, you simply need to print out fields `$2` and `$4` to catch the process ID and its process state. However, you only want to show zombie processes, so searching the output using the `grep` command will work. The command needed looks like the following:

```
$ ps -al | gawk '{print $2,$4}' | grep Z
Z 8779
```

This will work wonderfully for the snapshot report. Now, that you have all the pieces, let's put them together into a working script.

Creating the Snapshot Script

The first items you need to set up are the various variables to be used in the body of the script. These variables should be included at the top of the script to allow them to be easily modified in the future.

Our snapshot report will need a date stamp, so a date variable, `DATE`, will be included. For the performance statistics, you may have more than one disk to gather statistics for. Thus, an array variable, `DISKS_TO_MONITOR`, will include a list of all your system's disks. Because different

administrators set up different mail utilities, let's include a mail variable, `MAIL`, to designate what utility is locally used (see Chapter 25). Also, because a particular system's `mail` command may be located in a unique spot, you should use the `which` command (see Chapter 23) to find the `mail` command's location. Who to mail the report to is in the variable, `MAIL_TO`. And, finally, the name of the report is set in the variable, `REPORT`.

The variable section of the snapshot script looks like the following:

```
# Set Script Variables
#
DATE='date +%m%d%Y'
DISKS_TO_MONITOR="/dev/sda1 /dev/sda2"
MAIL='which mutt'
MAIL_TO=user
REPORT=/home/user/Documents/$DATE.rpt
```

Because you are generating a report and mailing it out of the same script, you need to save and later restore the file descriptor (see Chapter 14). The shell commands to save the file descriptor and redirect the output to `STDOUT` looks like this:

```
# Create Report File
#
exec 3>&1 #Save file descriptor
#
exec 1> $REPORT #direct output to rpt file.
```

The following code, prior to the `mail` command, will restore `STDOUT` to its saved file descriptor:

```
# Restore File Descriptor & Mail Report
#
exec 1>&3 #Restore output to STDOUT
```

With all the various pieces together, the script now looks

like the following:

```
#!/bin/bash
#
# Snapshot_Stats - produces a report for system stats
#####
# Set Script Variables
#
DATE=`date +%m%d%Y`
DISKS_TO_MONITOR="/dev/sda1 /dev/sda2"
MAIL='which mutt'
MAIL_TO=user
REPORT=/home/user/Documents/Snapshot_Stats_${DATE}.rpt
#
#####
# Create Report File
#
exec 3>&l #Save file descriptor
#
exec 1> $REPORT #direct output to rpt file.
#
#####
#
echo
echo -e "\t\tDaily System Report"
echo
#
#####
# Date Stamp the Report
#
echo -e "Today is " `date +%m/%d/%Y`
echo
#
#####
# 1) Gather System Uptime Statistics
#
echo -e "System has been \c"
uptime | sed -n '/,/s/,/ /gp' |
    gawk '{if ($4 == "days" || $4 == "day")
        {print $2,$3,$4,$5}
        else {print $2,$3}}'
#
#####
# 2) Gather Disk Usage Statistics
#
echo
for DISK in $DISKS_TO_MONITOR      #loop to check disk space
do
    echo -e "$DISK usage: \c"
    df -h $DISK | sed -n '/% \//p' | gawk '{print $5}'
```

```

done
#
#####
# 3) Gather Memory Usage Statistics
#
echo
echo -e "Memory Usage: \c"
#
free | sed -n '2p' |
    gawk 'x = int(($3 / $2) *100) {print x}' |
    sed 's/%/\'
#
#####
# 4) Gather Number of Zombie Processes
#
echo
ZOMBIE_CHECK='ps -al | gawk '{print $2,$4}' | grep Z'
#
if [ "$ZOMBIE_CHECK" = "" ]
then
    echo "No Zombie Process on System at this Time"
else
    echo "Current System Zombie Processes"
    ps -al | gawk '{print $2,$4}' | grep Z
fi
echo
#####
# Restore File Descriptor & Mail Report
#
exec 1>&3      #Restore output to STDOUT
#
$MAIL -a $REPORT -s "System Statistics Report for $DATE"
-- $MAIL_TO < /dev/null
#
#####
# Clean up
#
rm -f $REPORT
#

```

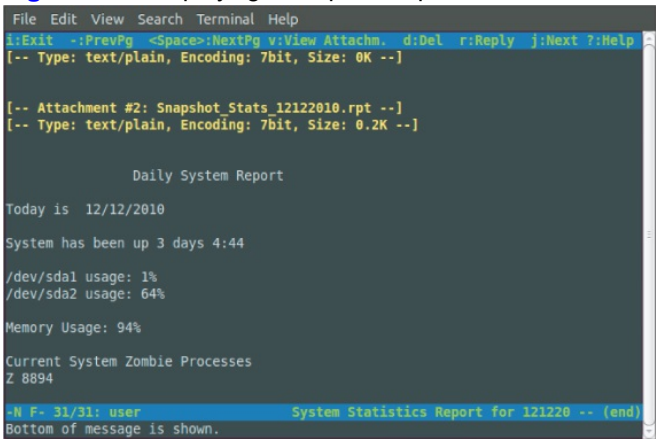
Notice that when the script reaches step #4, it first checks for zombie processes by using the results of the `ps` command put into the variable `ZOMBIE_CHECK`. If the variable is found to be empty, the script reports that no zombie processes exist. Otherwise, it lists out all current processes with a `Z` state.

Now, test out the new snapshot report script and see how well it performs.

```
$  
$ ./Snapshot_Stats
```

Not too much to see there. But if you get into your e-mail client, in this case Mutt, you will be able to see the snapshot report's output as shown in [Figure 27.1](#).

Figure 27.1 Displaying a snapshot report in mutt



The snapshot report is very useful for giving you a picture of a system's overall health. However, this type of report does not let you see various trends in performance and resource usage that occur over time. Therefore, let's move on to a set of scripts that will provide you with that type of

System Statistics Report

The core statistics of any Linux system are the CPU and memory usage. If these values start getting out of control, things can go wrong very quickly. This section demonstrates how to write scripts to help you monitor and track the CPU and memory usage on your Linux system over a period of time.

The Required Functions

The purpose of this first script will be to gather performance statistics into a data file for later use. Thus, you need to determine exactly what data you want to gather with your script and what commands to use.

For this script, we'll use the `uptime` command again. Except this time, instead of grabbing the amount of time the system has been up, we'll grab the number of users on the system. The number of users directly affects CPU and memory usage. Using `sed`, we'll search the `uptime` command's output text for the word "users." After the word is located in the current line, have `sed` just remove the rest of the text line, including the word "users." This will leave the statistic you want, number of users on the system, as the last item in the line. To obtain that last data item, you can use the `gawk` variable, `NF`. The entire command line and its output will look like this:

```
uptime | sed 's/users.*$//' | gawk '{print $NF}'
4
```

System load is another good statistic you can parse out of the `uptime` command output. System load shows, on average, how busy a CPU is on a system. A load average of 1 means a single CPU is busy all the time. However, if you have two CPUs on your system, this means that each one is only 50 percent busy. System load from the `uptime` command comes in averages for the last minute, for the last five minutes, and for the last fifteen minutes. We'll get the last fifteen minutes' average system load statistic, which is conveniently located at the end of the output text line:

```
$ uptime | gawk '{print $NF}'
0.19
```

Another great command for extracting system information is `vmstat`. Here's an example of the output from the `vmstat` command:

```
$ vmstat
procs -----memory----- ---swap-- -----io----- -system-- ----cpu-----
r b   swpd   free   buff   cache   si so       bi bo       in cs   us sy id wa
0 0    11328  42608  165444  502500   0 0         1  2        68 10    1  0 99  0
```

The first time you run the `vmstat` command, it displays the average values since the last reboot. To get the current statistics, you must run the `vmstat` command with command line parameters:

```
$ vmstat 1 2
procs -----memory----- ---swap-- -----io----- -system-- ----cpu-----
r b   swpd   free   buff   cache   si so       bi bo       in cs   us sy id wa
0 0    11328  43112  165452  502444   0 0         1  2        68 10    1  0 99  0
0 0    11328  40540  165452  505064   0 0         0  0        58 177    1  1 98  0
```

The second line contains the current statistics for the Linux system. As you can see, the output from the `vmstat`

command is somewhat cryptic. [Table 27.1](#) explains what each of the symbols means.

Table 27.1 The vmstat Output Symbols

Symbol	Description
r	The number of processes waiting for CPU time
b	The number of processes in uninterruptible sleep
swpd	The amount of virtual memory used (in MB)
free	The amount of physical memory not used (in MB)
buff	The amount of memory used as buffer space (in MB)
cache	The amount of memory used as cache space (in MB)
si	The amount of memory swapped in from disk (in MB)
so	The amount of memory swapped out to disk (in MB)
bi	Number of blocks received from a block device
bo	Number of blocks sent to a block device
in	The number of CPU interrupts per second
cs	The number of CPU context switches per second
us	Percent of CPU time spent running non-kernel code
sy	Percent of CPU time spent running kernel code
id	Percent of CPU time spent idle
wa	Percent of CPU time spent waiting for I/O

This is a lot of information. The free memory and percent of CPU time spent idle statistics will work for our purposes here.

You may have also noticed that the output from the `vmstat` command includes table heading information, which you obviously don't want in your data. To solve that problem,

use the `sed` command to display only lines that have a numerical value in them:

```
$ vmstat 1 2 | sed -n '/[0-9]/p'
1 0 11328 38524 165476 506548 0 0 1 2 68 10 1 0 99 0
0 0 11328 35820 165476 509528 0 0 0 0 82 160 1 1 98 0
```

That's better, but now you need to get only the second line of data. Another call to `sed` can solve that problem:

```
$ vmstat 1 2 | sed -n '/[0-9]/p' | sed -n '2p'
0 0 11328 36060 165484 509560 0 0 0 0 58 175 1 1 99 0
```

Now you can easily extract the data value you want using `gawk`.

Finally, you'll want to tag each performance data record with a date and timestamp to indicate when the statistics were taken. Simply use the `date` command, specifying the format you want in the data record.

```
$ date +"%m/%d/%Y %k:%M:%S"
12/12/2010 13:55:31
```

Now consider how you want to record the performance data values. For data that you sample on a regular basis, often it's best to output the data directly to a log file. You can create this file in your `$HOME` directory, appending data each time you run the shell script. When you want to see the results, you simply view the log file.

You'll want to ensure that the data in the log file can be read easily. There are many different methods you can use to format the data in the log file. A popular format is comma-separated values (CSV). This format places each record of data on a separate line, and separates the data fields in the record with commas. This is a popular format

as it's easily imported into spreadsheets, databases, and reporting utilities.

The script will save data in a CSV-formatted file. The reporting script we create later on in this chapter will use this file.

Creating the Capture Script

This script, used to capture the data, is really pretty simple once you have all the needed functions in place. Here's what it looks like:

```
#!/bin/bash
#
# Capture_Stats - Gather System Performance Statistics
#####
# Set Script Variables
#
REPORT_FILE=/home/user/Documents/capstats.csv
DATE='date +%m/%d/%Y'
TIME='date +%k:%M:%S'
#
#####
# Gather Performance Statistics
#
USERS='uptime | sed 's/users.*$/ /' | gawk '{print $NF}'
LOAD='uptime | gawk '{print $NF}'
#
FREE='vmstat 1 2 | sed -n '/[0-9]/p' | sed -n '2p' |
gawk '{print $4}'
IDLE='vmstat 1 2 | sed -n '/[0-9]/p' | sed -n '2p' |
gawk '{print $15}'
#
#####
# Send Statistics to Report File
#
echo "$DATE,$TIME,$USERS,$LOAD,$FREE,$IDLE" >> $REPORT_FILE
#
```

This script parses out the statistics from the `uptime` and `vmstat` commands and places them into variables. Those

variables are then written to the log file, `REPORT_FILE`, separated by commas, along with a date and time stamp. Notice that the data is appended to the log file using the `>>` redirection symbol. This will allow you to continue adding data onto the log file for as long as you deem necessary.

After creating the script, you should probably test it from the command line before having it run regularly from your cron table:

```
$ cat capstats.csv
12/09/2010, 9:06:50,2,0.29,645988,99
12/09/2010, 9:07:55,2,0.28,620252,100
12/09/2010, 9:40:51,3,0.37,474740,100
12/10/2010,14:36:46,3,0.30,46640,98
12/12/2010, 7:16:26,4,0.25,27308,98
12/12/2010,13:28:53,4,0.42,58832,100
```

The next script will create a report from the CSV data and e-mail it to the appropriate individual(s). Because the reporting script will be a separate script, you can have cron run it on a different time schedule than the `Capture_Stats` script. This will give you flexibility in how much data you capture before it is reported.

Generating the Report Script

Now that you have a file full of raw data, you can start working on the script to generate a nice report. The best tool for this is the `gawk` command.

The `gawk` command will allow you to extract the raw data from the CSV file and present it in any manner desired. First, test this from the command line, using the new `capstats.csv` file created by the `Capture_Stats` script:

```
cat capstats.csv |
> gawk -F, '{printf "%s %s - %s\n", $1, $2, $4}'
12/09/2010 9:06:50 - 0.29
12/09/2010 9:07:55 - 0.28
12/09/2010 9:40:51 - 0.37
12/10/2010 14:36:46 - 0.30
12/12/2010 7:16:26 - 0.25
12/12/2010 13:28:53 - 0.42
```

You need to use the `-F` option for the `gawk` command to define the comma as the field separator character in your data. After that, you can retrieve each individual data field and display it as you need with the `printf` feature.

For the report, we'll be using HTML format. HTML has been the standard method for formatting Web pages for years. It uses simple tags to delineate data types within the Web page. However, HTML is not just for Web pages. You'll often find HTML used in e-mail messages as well. Depending on your mail utility (see Chapter 25), you may or may not be able to view an embedded HTML e-mail document. A better solution is to create the HTML report and attach it to the e-mail message.

Creating the report with HTML allows you to produce a nicely formatted report with a minimum amount of work. The program that displays the report will do all the hard work of formatting and displaying the report. All you need to do is insert the appropriate HTML tags to format the data.

The easiest way to display spreadsheet-style data in HTML is with the `<table>` tag. The table tag allows you to create a table with rows and cells, which are called "divisions" in HTML. You define the start of a row using the `<tr>` tag, and the end of the row with the `</tr>` tag. Similarly,

you define cells using the `<td>` and `</td>` tag pair.

The HTML for a full table looks like this:

```
<html>
<body>
<h2>Report title</h2>
<table border="1">
<tr>
  <td>Date</td><td>Time</td><td>Users</td>
  <td>Load</td><td>Free Memory</td><td>%CPU Idle</td>
</tr>
<tr>
  <td>12/09/2010</td><td>11:00:00</td><td>4</td>
  <td>0.26</td><td>57076</td><td>87</td>
</tr>
</table>
</body>
</html>
```

Each data record is part of a `<tr>...</tr>` tag pair. Each data field is within its own `<td>...</td>` tag pair.

When you display the HTML report in a browser, the browser creates the table automatically for you, as shown in [Figure 27.2](#).

[Figure 27.2](#) Displaying data in an HTML table

File Edit View History Bookmarks Tools Help

file:///tmp/capstats.html

Most Visited Getting Started Latest Headlines

File:///tmp/capstats.html

Report for 12/12/2010

Date	Time	Users	Load	Free Memory	%CPU Idle
12/09/2010	9:06:50	2	0.29	645988	99
12/09/2010	9:07:55	2	0.28	620252	100
12/09/2010	9:40:51	3	0.37	474740	100
12/10/2010	14:36:46	3	0.30	46640	98
12/12/2010	7:16:26	4	0.25	27308	98
12/12/2010	13:28:53	4	0.42	58832	100

For the script, all you need to do is generate the HTML heading code by using `echo` commands, generate the data HTML code by using the `gawk` command, and then close out the table, again by using `echo` commands.

Here's the `Report_Stats` script, which will generate the performance report and mail it off:

```
#!/bin/bash
#
# Report_Stats - Generates Rpt from Captured Perf Stats
#####
# Set Script Variables
#
REPORT_FILE=/home/user/Documents/capstats.csv
TEMP_FILE=/home/user/Documents/capstats.html
#
DATE='date +%m/%d/%Y'
```

```
MAIL='which mutt'
MAIL_TO=user
#
#####
# Create Report Header
#
echo "<html><body><h2>Report for $DATE</h2>" > $TEMP_FILE
echo "<table border='1\">" >> $TEMP_FILE
echo "<tr><td>Date</td><td>Time</td><td>Users</td>" >> $TEMP_FILE
echo "<td>Load</td><td>Free Memory</td><td>%CPU Idle</td></tr>" >>
$TEMP_FILE
#
#####
# Place Performance Stats in Report
#
cat $REPORT_FILE | gawk -F, '{
printf "<tr><td>%s</td><td>%s</td><td>%s</td>", $1, $2, $3;
printf "<td>%s</td><td>%s</td><td>%s</td>\n</tr>\n", $4, $5, $6;
}' >> $TEMP_FILE
#
echo "</table></body></html>" >> $TEMP_FILE
#
#####
# Mail Performance Report & Clean up
#
$MAIL -a $TEMP_FILE -s "Performance Report $DATE"
-- $MAIL_TO < /dev/null
#
rm -f $TEMP_FILE
#
```

We use `mutt` to mail off our performance report, but you can change it to your local e-mail client. The name of the temporary HTML report file can also be modified to better suit your needs.

Caution

Most e-mail clients can automatically detect the type of file attachment by the file extension. For this reason, you should end your report file name with `.html`.

Running the Script

After creating the `Report_Stats` script, give it a test run from the command line and see what happens:

```
$  
$ ./Report_Stats
```

Well, that wasn't too exciting. The real test now is to view your mail message, preferably in a graphical e-mail client such as KMail or Evolution. [Figure 27.3](#) demonstrates how to view the message from the Evolution e-mail client. Notice how the data is all nicely formatted using the HTML tables, just as it was when viewing from the browser!

[Figure 27.3](#) Viewing the report attachment in Evolution

► 1 Attachment Save As

HTML document attachment (capstats.html)

Date	Time	Users	Load	Free Memory	%CPU Idle
12/09/2010	9:06:50	2	0.29	645988	99
12/09/2010	9:07:55	2	0.28	620252	100
12/09/2010	9:40:51	3	0.37	474740	100
12/10/2010	14:36:46	3	0.30	46640	98
12/12/2010	7:16:26	4	0.25	27308	98
12/12/2010	13:28:53	4	0.42	58832	100

Now you have scripts to generate two different reports that will help you track system performance and avoid problems due to a lack of resources, and so on. However, even excellent system monitoring cannot stop unexpected problems from occurring. Quickly solving those unexpected problems and learning from them can be helpful in minimizing their effects in the future. Therefore, in the next section, we'll take you through creating advanced database scripts that will allow you to track problems.

Problem-Tracking Database

Every system is doomed to experience the unexpected. However, analysis of a problem will allow you to make changes that will minimize future unexpected difficulties. It is also helpful to have a record of how you solved a problem, just in case it does happen again. There is an even bigger benefit when you work with a team of individuals who can record problems/solutions they experience. The problem/solution information can be more easily disseminated amongst the team members.

In this section, we take you through the process of setting up a basic problem-tracking database. Also, we help you create a few advanced scripts to record, update, and review the problems and their solutions.

Creating the Database

Planning the database structure will be the most important part of creating the problem-tracking database. You need to determine what kind of information you would like to track and who will have access to modify and use the data stored within the database.

The primary purpose for this database is to record a problem and how it was fixed so that you can use the information to solve future similar problems more quickly. To achieve this purpose, you should track the following information:

- The date the problem was reported
- The date the problem was fixed
- A description of the problem and its symptoms
- A description of the problem's resolution

We'll use `mysql` (see Chapter 23) as the database system for our problem-tracking database. To begin, you create an empty database:

```
$ $ mysql -u root -p
Enter password:
Welcome to the MySQL monitor. Commands end with ; or \g.
...
mysql> CREATE DATABASE Problem_Trek;
Query OK, 1 row affected (0.02 sec)

mysql>
```

With an empty database created, a table can now be added. We'll create one table, `problem_logger`, with four fields to contain the problem information. In addition, you will need one more field, `id_number`, to act as a primary key for your database. [Table 27.2](#) shows the five fields needed.

[Table 27.2](#) The Problem_Logger Table Fields

Field Name	Data Type
<code>id_number</code>	integer
<code>report_date</code>	date
<code>fixed_date</code>	date
<code>prob_symptoms</code>	text
<code>prob_solutions</code>	text

Creating this table in `mysql` is fairly straightforward:

```
mysql> USE Problem_Trek;
Database changed
```

```
mysql>
mysql> CREATE TABLE problem_logger (
-> id_number int not null,
-> report_date Date,
-> fixed_date Date,
-> prob_symptoms text,
-> prob_solutions text,
-> primary key (id_number));
Query OK, 0 rows affected (0.04 sec)
```

```
mysql>
```

Note that the `USE` command was issued to ensure we are creating our new table in the correct database. Let's take another look at the new `problem_logger` table using the `DESCRIBE` command. This is a nice clear way to display the fields added to our table:

```
mysql> DESCRIBE problem_logger;
+-----+-----+-----+-----+-----+-----+
| Field          | Type   | Null | Key | Default | Extra |
+-----+-----+-----+-----+-----+-----+
| id_number      | int(11)| NO   | PRI | NULL    |       |
| report_date    | date   | YES  |     | NULL    |       |
| fixed_date     | date   | YES  |     | NULL    |       |
| prob_symptoms  | text   | YES  |     | NULL    |       |
| prob_solutions | text   | YES  |     | NULL    |       |
+-----+-----+-----+-----+-----+-----+
5 rows in set (0.01 sec)
```

```
mysql>
```

So far, we had to use the `root` account to create the database and the table. Now that those steps are completed, the `root` account is no longer necessary for your work here. As you learned in Chapter 23, it is not a good practice to use the `root` account for database scripts. Thus, you need to add a new account to the database with the appropriate permissions. This account will be the one the scripts use to utilize the `Problem_Trek` database:

```
mysql> GRANT SELECT,INSERT,DELETE,UPDATE ON Problem_Trek.* TO
```



```
> cbres IDENTIFIED BY 'test_password';
Query OK, 0 rows affected (0.03 sec)
```

```
mysql>
```

Once you have added a new user account to the database, be sure to test it out:

```
$ mysql Problem_Trek -u cbres -p
Enter password:
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Welcome to the MySQL monitor. Commands end with ; or \g.
...
mysql>
```

The new account works perfectly! Now, one more step and you will be ready to start working on the needed functions for our database scripts. Using your favorite text editor, create the special configuration file, `$HOME/.my.cnf`, containing the password you set for the new user account created earlier. This will allow the scripts to access the `mysql` program without having to include a password within them.

```
$ cat $HOME/.my.cnf
[client]
password=test_password
$
```

And don't forget to restrict it to your access only:

```
$ chmod 400 $HOME/.my.cnf
```

Now the `Problem_Trek` database has been created and appropriate access has been set up for it. You are ready to start working on the database scripts.

Recording a Problem

The first step in tracking a problem is recording its

symptoms. Because recording a problem can be time-consuming using `mysql` commands, we will simplify the process by wrapping an easy-to-use script around them.

The Required Functions

To put problem information into your `problem_logger` table, you can use the `mysql INSERT` command (see Chapter 23). To insert a new record into the `problem_logger` table, you issue these commands:

```
$
$ mysql Problem_Trek -u cbres -p
Enter password:
...
mysql> INSERT INTO problem_logger VALUES (
-> 1012111322,
-> 20101211,
-> 0,
-> "When trying to run script Capture_Stats, getting message: bash:
./Capture_Stats: Permission denied",
-> "");
Query OK, 1 row affected (0.02 sec)

mysql>
```

That was pretty easy. You can look at the record you just put into the table by using the `SELECT` command. By specifying the new record's ID number to `WHERE`, only the record just entered will be shown:

```
mysql> SELECT * FROM problem_logger WHERE id_number=1012111322;
+-----+-----+-----+-----+-----+
| id_number | report_date | fixed_date | prob_symptoms | prob_solutions |
+-----+-----+-----+-----+-----+
| 1012111322 | 2010-12-11 | 0000-00-00 | When trying to run script Capture_Stats, getting message: bash:./Capture_Stats: Permission denied |
+-----+-----+-----+-----+-----+
1 row in set (0.00 sec)

mysql>
```

While you can see the record entered, it is hard to read the data in this format. To make the record easier to read, you can add a `\G` on the end of the command instead of a semicolon (`;`):

```
mysql> SELECT * FROM problem_logger WHERE id_number=1012111322\G
***** 1. row *****
id_number: 1012111322
report_date: 2010-12-11
fixed_date: 0000-00-00
prob_symptoms: When trying to run script Capture_Stats,
  getting message: bash:./Capture_Stats: Permission denied
prob_solutions:
1 row in set (0.00 sec)

mysql>
```

That looks much better and will work well within a script. Those two `mysql` commands are the main functions needed for the script. Thus, we can begin to build the `Record_Problem` script.

Creating the Script

In Chapter 23, you learned that it was wise to use the `which` command to find the `mysql` command on your particular system and then assign it to an environment variable. You can take that a step further and add the database name and user account for `mysql` to use:

```
#
MYSQL='which mysql'" Problem_Trek -u cbres"
#
```

The `Record_Problem` is an interactive script. To keep it simple, you can ask questions to obtain the needed information. As the script obtains answers, it will assign them to various variables. Once all the data is gathered,

you use the variable names along with the `INSERT` command to enter the information into the database. This way, the script user can be blissfully ignorant of the necessary `mysql` commands:

```
INSERT INTO problem_logger VALUES (  
    $ID_NUMBER,  
    $REPORT_DATE,  
    $FIXED_DATE,  
    "$PROB_SYMPTOMS",  
    "$PROB_SOLUTIONS");
```

The `ID_NUMBER` variable will be used for the primary key field in the table. Because it is a primary key field, the ID numbers will need to be unique for each problem recorded. To keep it simple for the script user, an automatically generated ID number should be used. To keep it unique, you can use a combination of the current date and time to produce the necessary 10-digit key:

```
#  
ID_NUMBER=`date +%y%m%d%H%M`  
#
```

Because this script is only for recording the problem, you don't necessarily have data yet for the `FIXED_DATE` and `PROB_SOLUTIONS` variables. Thus, set both of these variables to the appropriate null values before inserting the record into the table:

```
#  
# Set Fixed Date & Problem Solution to null for now  
FIXED_DATE=0  
PROB_SOLUTIONS=""  
#
```

To show the successful addition of the problem record, add a `SELECT` statement to the script. To keep the entire database's contents from displaying, the `WHERE` command

will use the ID_NUMBER variable to select only the record just added:

```
SELECT * FROM problem_logger WHERE id_number=$ID_NUMBER\G
```

Putting all these pieces together will result in a script as follows:

```
#!/bin/bash
#
# Record_Problem - records system problems in database
#####
# Determine mysql location & put into variable
#
MYSQL='which mysql'" Problem_Trek -u cbres"
#
#####
# Create Record Id & Report_Date
#
ID_NUMBER='date +%y%m%d%H%M'
#
REPORT_DATE='date +%Y%m%d'
#
#####
# Acquire information to put into table
#
echo
echo -e "Briefly describe the problem & its symptoms: \c"
#
read ANSWER
PROB_SYMPTOMS=$ANSWER
#
# Set Fixed Date & Problem Solution to null for now
FIXED_DATE=0
PROB_SOLUTIONS=""
#
#####
# Insert acquired information into table
#
#
echo
echo "Problem recorded as follows:"
echo
$MYSQL <<EOF
INSERT INTO problem_logger VALUES (
    $ID_NUMBER,
    $REPORT_DATE,
    $FIXED_DATE,
    "$PROB SYMPTOMS".
```

```

        "$PROB_SOLUTIONS");
SELECT * FROM problem_logger WHERE id_number=$ID_NUMBER\G
EOF
#
#####

```

Now, test out the `Record_Problem` script and see if it's ready for use:

```
$ ./Record_Problem
```

```
Briefly describe the problem & its symptoms:
```

```
Running yum to install software and
getting 'not found' message.
```

```
Problem recorded as follows:
```

```

***** 1. row *****
id_number: 1012111510
report_date: 2010-12-11
fixed_date: 0000-00-00
prob_symptoms:
Running yum to install software and getting 'not found' message.
prob_solutions:
$

```

That works wonderfully! You can see how it is so much easier to use a script like this to record problems than trying to insert them using the individual `mysql` commands.

Now that you have a way to record the problems, you need a way to update the records. Let's move on to creating the `Update_Problem` script.

Updating a Problem

The next step in tracking a problem is recording the problem's solution. The script created in this section will allow that information to be added to the problem's record.

The Required Functions

Updating records using `mysql` is very straightforward. Once all the necessary data is obtained, you just use the `mysql UPDATE` command:

```
UPDATE problem_logger SET
    prob_solutions="$PROB_SOLUTIONS",
    fixed_date=$FIXED_DATE
WHERE id_number=$ID_NUMBER;
```

That's pretty easy. However, putting together the `Update_Problem` script will take a little more than just using the `UPDATE` function.

Creating the Script

As in the previous script, we want to keep things fairly easy for the script user. Several additions are needed to reach that goal.

The primary key for the `Problem_Trek` database, `id_number`, is needed to identify the record to update. If the ID number is known, it can be passed as a parameter to the script (see Chapter 13). All you need to do is check for the parameter and assign it to the variable `ID_NUMBER`. If the ID number is not passed as a parameter, the script will ask for it:

```
if [ $# -eq 0 ] #Check if id number was passed
then          #if not passed ask for it.
#
...
    echo
    echo "What is the ID number for the"
    echo -e "problem you want to update?: \c"
    read ANSWER
    ID_NUMBER=$ANSWER
else
    ID_NUMBER=$1
```

However, asking for the problem ID number is a little unfair. Not too many individuals carry around in their heads a 10-digit number for identifying a particular problem they are solving. Thus, you should provide a list of the various unresolved problems in the database, making it easier on the script user.

First, check if there are any unresolved records. The fields that will be empty are `fixed_date` and `prob_solutions`. If either of these fields is empty, the record needs to be updated; with `SELECT` you can use the `OR` command in the `WHERE` statement to search both fields. The code to accomplish this looks like the following:

```
RECORDS_EXIST='$MYSQL -Bse `SELECT id_number FROM problem_logger  
WHERE fixed_date="0000-00-00" OR prob_solutions=""`'
```

Notice that the output of the command is assigned to the variable `RECORDS_EXIST`. The `-Bse` options had to be used on the `mysql` command in order to do this. They allow the command to be executed in batch (all at once) and then exit. If the `RECORDS_EXIST` variable contains data, then there are records that need updating. Display them to the screen using the following:

```
#  
  
if [ "$RECORDS_EXIST" != "" ]  
then  
echo  
echo "The following record(s) need updating..."  
$MYSQL <<EOF  
SELECT id_number, report_date, prob_symptoms  
FROM problem_logger  
WHERE fixed_date="0000-00-00" OR  
prob_solutions=""\G  
  
EOF
```


#

This will provide a nice list of unfinished problem records along with their ID numbers. Picking an ID number from a list is much easier than trying to remember the numbers!

Now you can put together an entire script for updating records that will be simple to use. It should look similar to the following:

```
#!/bin/bash
#
# Update_Problem - updates problem record in database
#####
# Determine sql location & set variable
#
MYSQL='which mysql' Problem_Trek -u cbres"
#
#####
# Obtain Record Id
#
if [ $# -eq 0 ] #Check if id number was passed
then          #if not passed ask for it.
#
#       Check if any unfinished records exist.
RECORDS_EXIST=`$MYSQL -Bse 'SELECT id_number FROM problem_logger
WHERE fixed_date="0000-00-00" OR prob_solutions=""'`
#
if [ "$RECORDS_EXIST" != "" ]
then
echo
echo "The following record(s) need updating..."
$MYSQL <<EOF
SELECT id_number, report_date, prob_symptoms
FROM problem_logger
WHERE fixed_date="0000-00-00" OR
      prob_solutions=""\G
EOF
fi
#
echo
echo "What is the ID number for the"
echo -e "problem you want to update?: \c"
read ANSWER
```

```

ID_NUMBER=$ANSWER
else
    ID_NUMBER=$1
fi
#
#####
# Obtain Solution (aka Fixed) Date
#
echo
echo -e "Was problem solved today? (y/n) \c"
read ANSWER
#
case $ANSWER in
y|Y|YES|yes|Yes|yEs|yeS|YES|yES )
#
    FIXED_DATE='date +%Y%m%d'
;;
*)
# If answer is anything but yes, ask for date
echo
echo -e "What was the date of resolution? [YYYYMMDD] \c"
read ANSWER
#
    FIXED_DATE=$ANSWER
;;
esac
#
#####
# Acquire problem solution
#
echo
echo -e "Briefly describe the problem solution: \c"
#
read ANSWER
PROB_SOLUTIONS=$ANSWER
#
#####
# Update problem record
#
#
echo
echo "Problem record updated as follows:"
echo
$MYSQL <<EOF
UPDATE problem_logger SET
    prob_solutions="$PROB_SOLUTIONS",
    fixed_date=$FIXED_DATE
WHERE id_number=$ID_NUMBER;
SELECT * FROM problem_logger WHERE id_number=$ID_NUMBER\G
EOF

```

Give this script a test and see how well it performs.

```
$ ./Update_Problem

The following record(s) need updating...
***** 1. row *****
id_number: 1012111624
report_date: 2010-12-11
prob_symptoms: Running yum to install software and getting
'not found' message.

What is the ID number for the
problem you want to update?: 1012111624

Was problem solved today? (y/n) y

Briefly describe the problem solution: Network service was down.
Issued the command: service network restart

Problem record updated as follows:

***** 1. row *****
id_number: 1012111624

report_date: 2010-12-11
fixed_date: 2010-12-11
prob_symptoms: Running yum to install software and getting
'not found' message.
prob_solutions: Network service was down.
Issued the command: service network restart
$
```

This works very well and is rather simple to use.

An Addition

Now to make this process even easier for the script user, let's modify the `Record_Problem` script to call `Update_Problem`, if desired. Sometimes, when you get busy, it is easier to record problems in “batch.” At the very end of the original `Record_Problem` script, you can add the following code:

```
# Check if want to enter a solution now
```

```

echo
echo -e "Do you have a solution yet? (y/n) \c"
read ANSWER
#
case $ANSWER in
y|Y|YES|yes|Yes|yEs|yeS|YES|yES )
    ./$HOME/scripts/Update_Problem $ID_NUMBER
#
;;
*)
# If answer is anything but yes, just exit script
;;
esac
#####

```

Now let's test this modification to see if it works:

```
$ ./Record_Problem
```

Briefly describe the problem & its symptoms:

Moved script from Fedora to Ubuntu and was not working properly.

Problem recorded as follows:

```

***** 1. row *****
id_number: 1012111631
report_date: 2010-12-11
fixed_date: 0000-00-00
prob_symptoms:
Moved script from Fedora to Ubuntu and was not working properly.
prob_solutions:
Do you have a solution yet? (y/n) y

Was problem solved today? (y/n) y

```

Briefly describe the problem solution:

Added #!/bin/bash to the top of the script.

This is needed because Ubuntu's default shell is dash.

Problem record updated as follows:

```

***** 1. row *****
id_number: 1012111631
report_date: 2010-12-11
fixed_date: 2010-12-11
prob_symptoms:
Moved script from Fedora to Ubuntu and was not working properly.
prob_solutions: Added #!/bin/bash to the top of the script.
This is needed because Ubuntu's default shell is dash.

```

This works perfectly. Now you have made it very easy to record and update the records in the `Problem_Trek` database.

Finding a Problem

The real beauty of having a database such as `Problem_Trek` is the ability to review the records. Looking at past problems and their solutions will allow you to resolve repeat problems quickly and spot trends.

The Required Functions

To find records in the table to review, you can again use the `SELECT` command. We have used it previously in this chapter for purposes such as viewing newly added records in the `problem_logger` table:

```
SELECT * FROM problem_logger WHERE id_number=$ID_NUMBER\G
```

For the `Find_Problem` script, a slight change needs to be made to the `SELECT` command we used previously. So that old problems can be easily found, the script will ask for a keyword to be entered. Thus, if you remembered that you had a problem with `yum` last month, you would enter the keyword “yum.” To search the database for a record containing that keyword, add on the command `LIKE` to the script's `SELECT` statement. The code looks similar to the following:

```
mysql> SELECT * FROM problem_logger
-> WHERE prob_symptoms LIKE 'yum'\G
Empty set (0.00 sec)
```

Wait, that didn't work! This is due to the fact that wildcards are needed. To find the word “yum” buried in the `prob_symptoms` field, a wildcard needs to be added to both sides of the keyword. In effect, you are telling `mysql` to find a record whose `prob_symptoms` field contains the word “yum.” The wildcard symbol to use is the percentage sign (%). Thus, your code should look like this:

```
mysql> SELECT * FROM problem_logger
-> WHERE prob_symptoms LIKE '%yum%'\G
***** 1. row *****
id_number: 1012111624
report_date: 2010-12-11
fixed_date: 2010-12-11
prob_symptoms: Running yum to install software and getting
'not found' message.
prob_solutions: Network service was down.
Issued the command: service network restart
1 row in set (0.00 sec)
```

You will want the `SELECT` command to search the `prob_solutions` field as well as the `prob_symptoms` field, which will allow a more thorough search for the desired record:

```
mysql> SELECT * FROM problem_logger WHERE
-> prob_symptoms LIKE '%yum%'
-> OR
-> prob_solutions LIKE '%yum%'\G
***** 1. row *****
id_number: 1012111624
...
```

Note

If you are worried about case sensitivity, don't be. By default, `mysql` will ignore case, so you can enter the keyword “yum” or the keyword “YUM”

To be an effective search through the `Problem_Trek`

database, more than one keyword should be allowed for a search. The `mysql` command `LIKE` cannot accommodate multiple keywords. However, you can use the regular expression command `REGEXP`. One of the nice things about `REGEXP` is that it will search the entire designated field for the keyword you specify—no wildcards needed. To search for multiple keywords, you put a logical `or`, a pipe symbol (`|`), between them. Thus, the code you need now looks like the following:

```
mysql> SELECT * FROM problem_logger WHERE
-> prob_symptoms REGEXP 'yum|dash'
-> OR
-> prob_solutions REGEXP 'yum|dash'\G
***** 1. row *****
id_number: 1012111624
report_date: 2010-12-11
fixed_date: 2010-12-11
prob_symptoms: Running yum to install software and getting
'not found' message.
prob_solutions: Network service was down.
Issued the command: service network restart
***** 2. row *****
id_number: 1012111631
report_date: 2010-12-11
fixed_date: 2010-12-11
prob_symptoms:
Moved script from Fedora to Ubuntu and was not working properly.
prob_solutions: Added #!/bin/bash to the top of the script.
This is needed because Ubuntu's default shell is dash.
2 rows in set (0.00 sec)
```

Now you have the anchor function needed for the `Find_Problem` script. Let's put the rest of the script together.

Creating the Script

It is best to allow keywords to be passed as parameters when the script is called. However, you have no idea how

many keywords may be passed. It could be five, it could be one, or it could be none! This is actually an easy dilemma to solve. Use an `if` test condition to see if any keywords were passed. If they were, you grab all of them at once using the `$@` parameter (see Chapter 13). If no keywords were passed, the script will need to ask for them. The code to accomplish this task looks like the following:

```
# Obtain Keyword(s)
#
if [ -n "$1" ]           #Check if a keyword was passed
then                     #Grab all the passed keywords
#
    KEYWORDS=$@          #Grab all the params as separate words, same string
#
else                     #Keyword(s) not passed. Ask for them.
    echo
    echo "What keywords would you like to search for?"
    echo -e "Please separate words by a space: \c"
    read ANSWER
    KEYWORDS=$ANSWER
fi
#
```

Once you have the keywords, the next step is to build the search statement. Remember that a logical `or (|)` will need to be placed between each pair of keywords in order for the `REGEXP` to search for them properly. This is easily accomplished with `sed`. Every space between in the `KEYWORDS` array variable can be replaced with a `|`.

```
KEYWORDS=`echo $KEYWORDS | sed 's/ /|/g'`
```

Now the script looks like this:

```
#!/bin/bash
#
# Find_Problem - finds problem records using keywords
#####
# Determine sql location & set variable
#
MYSQL=`which mysql` Problem_Trek -u cbres"
```



```
#####
# Obtain Keyword(s)
#
if [ -n "$1" ]          #Check if a keyword was passed
then                    #Grab all the passed keywords
#
    KEYWORDS=$@         #Grab all the params as separate words, same string
#
else                    #Keyword(s) not passed. Ask for them.
    echo
    echo "What keywords would you like to search for?"
    echo -e "Please separate words by a space: \c"
    read ANSWER
    KEYWORDS=$ANSWER
fi
#
#####
# Find problem record
#
echo
echo "The following was found using keywords: $KEYWORDS"
echo
#
KEYWORDS='echo $KEYWORDS | sed `s/ /\|/g`'
#
$MYSQL <<EOF
SELECT * FROM problem_logger WHERE
    prob_symptoms REGEXP '($KEYWORDS)'
    OR
    prob_solutions REGEXP '($KEYWORDS)'\G
EOF
#
```

And of course, the next step is to test your new script out. Because the multiple keyword issue has been in focus, let's first try the script passing multiple keywords at the command line:

```
$ ./Find_Problem yum dash
```

The following was found using keywords: yum dash

```
***** 1. row *****
id_number: 1012111624
...
prob_solutions: Added #!/bin/bash to the top of the script.
This is needed because Ubuntu's default shell is dash.
```

Now, let's test the script sending no keywords as parameters. Instead, we'll let the script ask us for the keywords.

```
$ ./Find_Problem

What keywords would you like to search for?
Please separate words by a space: Ubuntu

The following was found using keywords: Ubuntu

***** 1. row *****
id_number: 1012111631
...
prob_solutions: Added #!/bin/bash to the top of the script.
This is needed because Ubuntu's default shell is dash.
```

Both tests ran without any problems. Now finding problem records in the `Problem_Trek` database will be easy.

With all three of these scripts, `Record_Problem`, `Update_Problem`, and `Find_Problem`, you can start to create your own problem database. We hope the time you spend creating these scripts will be given back to you in the time you save solving and avoiding problems.

Summary

In this chapter, we reviewed some of the more advanced features you can put into shell scripts. Users frequently overlook the ways in which shell scripts can be useful on particular system administration tasks because the scripts need advanced features to operate properly.

The chapter demonstrated how to use shell performance utilities, such as `uptime`, `df`, and `free`, to gather a snapshot of

system performance and resource usage. The desired information was gleaned from the commands' output using `sed` and `gawk`, including some of their more advanced capabilities. After a report was created from the information, it was sent to a designated user account via an e-mail client.

The next section covered the creation of a more advanced performance monitor script using `uptime` and `vmstat`. Data was collected from these utilities and stored in a CSV file for use by other scripts and/or utilities. A second script was created to read the CSV file's performance data and create an HTML-formatted report. Because of its HTML format, the report could be read via an e-mail utility or within a favorite browser.

We ended the chapter by stepping through the creation of a problem-tracking database and the advanced shell scripts used to populate, update, and report on the database. We used the `mysql` database software and included commands such as `INSERT`, `UPDATE`, and `SELECT` in the scripts. We also covered difficult issues such as searching records with the `LIKE` command versus the regular expression command `REGEXP`.

All of the scripts written in this chapter are easily modified for use in your own unique environment. By pointing out some unique and advanced ways to put shell scripts to work, we hope we've made you think about some new ways to use shell scripts and new scripts you can create.

Thanks for joining us on this journey through the Linux command line and shell scripting. We hope that you've enjoyed the journey and have learned how to get around on the command line, and how to create shell scripts to save you time. But don't stop your command line education here. There's always something new being developed in the open source world, whether it's a new command line utility or a full-blown shell. Stay in touch with the Linux community and follow along with the new advances and features.

Appendix A

Quick Guide to bash Commands

In This Appendix

- The bash built-in commands
- The GNU additional shell commands
- The bash environment variables

As you've seen throughout this book, the bash shell contains lots of features and thus has lots of commands available. This appendix provides a concise guide to allow you to quickly look up a feature or command that you can use from the bash command line or from a bash shell script.

Built-in Commands

The bash shell includes many popular commands built into the shell. This provides for faster processing times when using these commands. [Table A.1](#) shows the built-in

commands available directly from the bash shell.

Table A.1 bash Built-in Commands

Command	Description
alias	Define an alias for the specified command.
bg	Resume a job in background mode.
bind	Bind a keyboard sequence to a readline function or macro.
break	Exit from a <code>for</code> , <code>while</code> , <code>select</code> , or <code>until</code> loop.
builtin	Execute the specified shell built-in command.
cd	Change the current directory to the specified directory.
caller	Return the context of any active subroutine call.
command	Execute the specified command without the normal shell lookup.
compgen	Generate possible completion matches for the specified word.
complete	Display how the specified words would be completed.
continue	Resume the next iteration of a <code>for</code> , <code>while</code> , <code>select</code> , or <code>until</code> loop.
declare	Declare a variable or variable type.
dirs	Display a list of currently remembered directories.
disown	Remove the specified jobs from the jobs table for the process.
echo	Display the specified string to STDOUT.
enable	Enable or disable the specified built-in shell command.
eval	Concatenate the specified arguments into a single command, and then execute the command.
exec	Replace the shell process with the specified command.
exit	Force the shell to exit with the specified exit status.
	Set the specified variables to be available for child shell

export	processes.
fc	Select a list of commands from the history list.
fg	Resume a job in foreground mode.
getopts	Parse the specified positional parameters.
hash	Find and remember the full pathname of the specified command.
help	Display a help file.
history	Display the command history.
jobs	List the active jobs.
kill	Send a system signal to the specified process ID (PID).
let	Evaluate each argument in a mathematical expression.
local	Create a limited-scope variable in a function.
logout	Exit a login shell.
popd	Remove entries from the directory stack.
printf	Display text using formatted strings.
pushd	Add a directory to the directory stack.
pwd	Display the pathname of the current working directory.
read	Read one line of data from STDIN and assign it to a variable.
readonly	Read one line of data from STDIN and assign it to a variable that can't be changed.
return	Force a function to exit with a value that can be retrieved by the calling script.
set	Set and display environment variable values and shell attributes.
shift	Rotate positional parameters down one position.
shopt	Toggle the values of variables controlling optional shell behavior.
	Suspend the execution of the shell until a SIGCONT signal is

suspend	received.
test	Return an exit status of 0 or 1 based on the specified condition.
times	Display the accumulated user and system
trap	Execute the specified command if the specified system signal is received.
type	Display how the specified word would be interpreted if used as a command.
ulimit	Set a limit on the specified resource for system users.
umask	Set default permissions for newly created files and directories.
unalias	Remove the specified alias.
unset	Remove the specified environment variable or shell attribute.
wait	Wait for the specified process to complete, and return the exit status.

The built-in commands provide higher performance than external commands, but the more built-in commands that are added to a shell, the more memory it consumes with commands that you may never use. The bash shell also contains external commands that provide extended functionality for the shell. These are discussed in the bash Commands section.

bash Commands

In addition to the built-in commands, the bash shell utilizes external commands to allow you to maneuver around the filesystem and manipulate files and directories. [Table A.2](#)

shows the common external commands you'll want to use when working in the bash shell.

Table A.2 The bash Shell External Commands

Command	Description
bzip2	Compression using the Burrows-Wheeler block sorting text compression algorithm and Huffman coding.
cat	List the contents of the specified file.
chage	Change the password expiration date for the specified system user account.
chfn	Change the specified user account's comment information.
chgrp	Change the default group of the specified file or directory.
chmod	Change system security permissions for the specified file or directory.
chown	Change the default owner of the specified file or directory.
chpasswd	Read a file of login name and password pairs and update the passwords.
chsh	Change the specified user account's default shell.
compress	Original Unix file compression utility.
cp	Copy the specified files to an alternate location.
date	Display the date in various formats.
df	Display current disk space statistics for all mounted devices.
du	Display disk usage statistics for the specified filepath.
file	View the file type of the specified file.
find	Perform a recursive search for files.
finger	Display information about user accounts on the Linux system or a remote system.
free	Check available and used memory on the system.

grep	Search a file for the specified text string.
groupadd	Create a new system group.
groupmod	Modify an existing system group.
gzip	The GNU Project's compression using Lempel-Ziv compression.
head	Display the first portion of the specified file's contents.
killall	Send a system signal to a running process based on process name.
less	Advanced viewing of file contents.
link	Create a link to a file using an alias name.
ls	List directory contents.
mkdir	Create the specified directory under the current directory.
more	List the contents of the specified file, pausing after each screen of data.
mount	Display or mount disk devices into the virtual filesystem.
mv	Rename a file.
nice	Run a command using a different priority level on the system.
passwd	Change the password for a system user account.
ps	Display information about the running processes on the system.
pwd	Display the current directory.
renice	Change the priority of a running application on the system.
rm	Delete the specified file.
rmdir	Delete the specified directory.
sort	Organize data in a data file based on the specified order.
stat	View the file statistics of the specified file.
sudo	Run an application as the root user account.

tail	Display the last portion of the specified file's contents.
tar	Archive data and directories into a single file.
touch	Create a new empty file, or update the timestamp on an existing file.
top	Display the active processes, showing vital system statistics.
umount	Remove a mounted disk device from the virtual filesystem.
uptime	Display information on how long the system has been running.
useradd	Create a new system user account.
userdel	Remove an existing system user account.
usermod	Modify an existing system user account.
vmstat	Produce a detailed report on memory and CPU usage on the system.
which	Find the location of an executable file
zip	Unix version of the Windows PKZIP program.

You can accomplish just about any task you need to on the command line using these commands.

Environment Variables

The bash shell also utilizes many environment variables. While environment variables aren't commands, they often affect how shell commands operate, so it's important to know the shell environment variables. [Table A.3](#) shows the default environment variables available in the bash shell.

[Table A.3](#) bash Shell Environment Variables

Variable	Description
\$*	Contains all of the command line parameters as a single text value.
\$@	Contains all of the command line parameters as separate text values.
\$#	The number of command line parameters.
\$?	The exit status of the most recently used foreground process.
\$-	The current command line option flags.
\$\$	The process ID (PID) of the current shell.
\$_	The PID of the most recently executed background process.
\$0	The name of the command from the command line.
\$_	The absolute pathname of the shell.
BASH	The full file name used to invoke the shell.
BASH_ARGC	The number of parameters in the current subroutine.
BASH_ARGV	An array containing all of the command line parameters specified.
BASH_COMMAND	The name of the command currently being executed.
BASH_ENV	When set, each bash script attempts to execute a startup file defined by this variable before running.
BASH_EXECUTION_STRING	The command used in the -c command line option.
BASH_LINENO	An array containing the line numbers of each command in the script.
BASH_REMATCH	An array containing text elements that match a specified regular expression.
BASH_SOURCE	An array containing source file names for the declared functions in the shell.
	The number of subshells spawned by the current

BASH_SUBSHELL	shell.
BASH_VERSION	The version number of the current instance of the bash shell.
BASH_VERSINFO	A variable array that contains the individual major and minor version numbers of the current instance of the bash shell.
COLUMNS	Contains the terminal width of the terminal used for the current instance of the bash shell.
COMP_CWORD	An index into the variable <code>COMP_WORDS</code> , which contains the current cursor position.
COMP_LINE	The current command line.
COMP_POINT	The index of the current cursor position relative to the beginning of the current command.
COM_WORDBREAKS	A set of characters used as word separators when performing word completion.
COMP_WORDS	A variable array that contains the individual words on the current command line.
COMPREPLY	A variable array that contains the possible completion codes generated by a shell function.
DIRSTACK	A variable array that contains the current contents of the directory stack.
EUID	The numeric effective user ID of the current user.
FCEDIT	The default editor used by the <code>fc</code> command.
FIGIGNORE	A colon-separated list of suffixes to ignore when performing file name completion.
FUNCNAME	The name of the currently executing shell function.
GLOBIGNORE	A colon-separated list of patterns defining the set of file names to be ignored by file name expansion.
GROUPS	A variable array containing the list of groups of which the current user is a member.

histchars	Up to three characters that control history expansion.
HISTCMD	The history number of the current command.
HISTCONTROL	Controls what commands are entered in the shell history list.
HISTFILE	The name of the file to save the shell history list (<code>.bash_history</code> by default).
HISTFILESIZE	The maximum number of lines to save in the history file.
HISTIGNORE	A colon-separated list of patterns used to decide which commands are ignored for the history file.
HISTSIZE	The maximum number of commands stored in the history file.
HOSTFILE	Contains the name of the file that should be read when the shell needs to complete a hostname.
HOSTNAME	The name of the current host.
HOSTTYPE	A string describing the machine the bash shell is running on.
IGNOREEOF	The number of consecutive EOF characters the shell must receive before exiting. If this value doesn't exist, the default is 1.
INPUTRC	The name of the readline initialization file (the default is <code>.inputrc</code>).
LANG	The locale category for the shell.
LC_ALL	Overrides the <code>LANG</code> variable, defining a locale category.
LC_COLLATE	Sets the collation order used when sorting string values.
LC_CTYPE	Determines the interpretation of characters used in file name expansion and pattern matching.
LC_MESSAGES	Determines the locale setting used when interpreting double-quoted strings preceded by a

	dollar sign.
LC_NUMERIC	Determines the locale setting used when formatting numbers.
LINENO	The line number in a script currently executing.
LINES	Defines the number of lines available on the terminal.
MACHTYPE	A string defining the system type in <i>cpu-company-system</i> format.
MAILCHECK	How often (in seconds) the shell should check for new mail (default is 60).
OLDPWD	The previous working directory used in the shell.
OPTERR	If set to 1, the bash shell displays errors generated by the <code>setopt</code> command.
OSTYPE	A string defining the operating system the shell is running on.
PIPESTATUS	A variable array containing a list of exit status values from the processes in the foreground process.
POSIXLY_CORRECT	If set, bash starts in POSIX mode.
PPID	The process ID (PID) of the bash shell's parent process.
PROMPT_COMMAND	If set, the command to execute before displaying the primary prompt.
PS1	The primary command line prompt string.
PS2	The secondary command line prompt string.
PS3	The prompt to use for the <code>select</code> command.
PS4	The prompt displayed before the command line is echoed if the bash <code>-x</code> parameter is used.
PWD	The current working directory.
RANDOM	Returns a random number between 0 and 32767. Assigning a value to this variable seeds the random

	number generator.
REPLY	The default variable for the <code>read</code> command.
SECONDS	The number of seconds since the shell was started. Assigning a value resets the timer to the value.
SHELLOPTS	A colon-separated list of enabled bash shell options.
SHLVL	Indicates the shell level, incremented by 1 each time a new bash shell is started.
TIMEFORMAT	A format specifying how the shell displays time values.
TMOUT	The value of how long (in seconds) the <code>select</code> and <code>read</code> commands should wait for input. The default of 0 indicates to wait indefinitely.
UID	The numeric real user ID of the current user.

You display the environment variables using the `set` built-in command. The default shell environment variables set at boot time can and often do vary between different Linux distributions.

Appendix B

Quick Guide to sed and gawk

In This Appendix

- The basics for using sed
- What you need to know about gawk

If you do any type of data handling in your shell scripts, most likely you'll need to use either the sed or gawk programs (and sometimes both). This appendix provides a quick reference for sed and gawk commands that come in handy when working with data in your shell scripts.

The sed Editor

The sed editor can manipulate data in a data stream based on commands you either enter into the command line or store in a command text file. It reads one line of data at a time from the input and matches that data with the supplied editor commands, changes data in the stream as specified in the commands, and then outputs the new data to `STDOUT`.

Starting the sed Editor

The format for using the sed command is:

```
sed options script file
```

The *options* parameters allow you to customize the behavior of the sed command and include the options shown in [Table B.1](#).

[Table B.1](#) The sed Command Options

Option	Description
<code>-e script</code>	Add commands specified in <i>script</i> to the commands run while processing the input.
<code>-f file</code>	Add the commands specified in the file <i>file</i> to the commands run while processing the input.
<code>-n</code>	Don't produce output for each command, but wait for the <code>print</code> command.

The *script* parameter specifies a single command to apply against the stream data. If more than one command is required, you must use either the `-e` option to specify them in the command line or the `-f` option to specify them in a separate file.

sed Commands

The sed editor script contains commands that sed processes for each line of data in the input stream. This section describes some of the more common sed commands you'll want to use.

Substitution

The `s` command substitutes text in the input stream. The format of the `s` command is:

```
s/pattern/replacement/flags
```

`pattern` is the text to replace, and `replacement` is the new text that sed will insert in its place.

The `flags` parameter controls how the substitution takes place. Four types of substitution flags are available:

- A number, indicating the pattern occurrence that should be replaced.
- `g`: Indicates that all occurrences of the text should be replaced.
- `p`: Indicates that the contents of the original line should be printed.
- `w file`: Indicates that the results of the substitution should be written to a file.

In the first type of substitution, you can specify which occurrence of the matching pattern the sed editor should replace. For example, you use a 2 to replace only the second occurrence of the pattern.

Addressing

By default, the commands you use in the sed editor apply to all lines of the text data. If you want to apply a command to only a specific line, or a group of lines, you must use *line addressing*.

There are two forms of line addressing in the sed editor:

- A numeric range of lines
- A text pattern that filters out a line

Both forms use the same format for specifying the address:

```
[address]command
```

When using numeric line addressing, you reference lines by their line position in the text stream. The sed editor assigns the first line in the text stream as line number one and continues sequentially for each new line.

```
$ sed '2,3s/dog/cat/' data1
```

The other method of restricting which lines a command applies to is a bit more complicated. The sed editor allows you to specify a text pattern that it uses to filter lines for the command. The format for this is:

```
/pattern/command
```

You must encapsulate the *pattern* you specify in forward slashes. The sed editor applies the command only to lines that contain the text pattern that you specify.

```
$ sed '/rich/s/bash/csh/' /etc/passwd
```

This filter finds the line that contains the text `rich`, and replaces the text `bash` with `csh`.

You can also group more than one command together for a specific address:

```
address {  
    command1  
    command2  
    command3}
```

The sed editor applies each of the commands you

specify only to lines that match the address specified. The sed editor will process each command listed on the address line(s):

```
$ sed '2{  
> s/fox/elephant/  
> s/dog/cat/  
> }' data1
```

The sed editor applies each of the substitutions to the second line in the data file.

Deleting Lines

The delete command, `d`, pretty much does what it says. It deletes any text lines that match the addressing scheme supplied. Be careful with the delete command, for if you forget to include an addressing scheme, all of the lines will be deleted from the stream:

```
$ sed 'd' data1
```

The delete command is obviously most useful when used in conjunction with a specified address. This allows you to delete specific lines of text from the data stream, either by line number:

```
$ sed '3d' data6
```

or a specific range of lines:

```
$ sed '2,3d' data6
```

The pattern-matching feature of the sed editor also applies to the delete command:

```
$ sed '/number 1/d' data6
```

Only lines matching the specified text are deleted from the

stream.

Inserting and Appending Text

As you would expect, like any other editor, the sed editor allows you to insert and append text lines to the data stream. The difference between the two actions can be confusing:

- The insert command (*i*) adds a new line before the specified line.
- The append command (*a*) adds a new line after the specified line.

The format of these two commands can be confusing: You can't use these commands on a single command line. You must specify the line to insert or append on a separate line by itself. The format for doing this is:

```
sed '[address]command\  
new line'
```

The text in *new line* appears in the sed editor output in the place you specify. Remember that when you use the insert command, the text appears before the data stream text:

```
$ echo "testing" | sed 'i\  
> This is a test'  
This is a test  
testing  
$
```

And when you use the append command, the text appears after the data stream text:

```
$ echo "testing" | sed 'a\  
> This is a test'  
testing  
This is a test
```

This allows you to insert text at the end of the normal text.

Changing Lines

The change command allows you to change the contents of an entire line of text in the data stream. It works the same as the insert and append commands, in that you must specify the new line separately from the rest of the `sed` command:

```
$ sed '3c\  
> This is a changed line of text.' data6
```

The backslash character is used to indicate the new line of data in the script.

Transform Command

The transform command (`y`) is the only `sed` editor command that operates on a single character. The transform command uses the format:

```
[address]y/inchars/outchars/
```

The transform command performs a one-to-one mapping of the `inchars` and the `outchars` values. The first character in `inchars` is converted to the first character in `outchars`. The second character in `inchars` is converted to the second character in `outchars`. This mapping continues throughout the length of the specified characters. If the `inchars` and `outchars` are not the same length, the `sed` editor will produce an error message.

Printing Lines

Similar to the `p` flag in the substitution command, the `p` command prints a line in the sed editor output. The most common use for the print command is for printing lines that contain matching text from a text pattern:

```
$ sed -n '/number 3/p' data6
This is line number 3.
$
```

The `print` command allows you to filter only specific lines of data from the input stream.

Writing to a File

The `w` command is used to write lines to a file. The format for the `w` command is:

```
[address]w filename
```

The *filename* can be specified as either a relative or absolute pathname, but in either case the person running the sed editor must have write permissions for the file. The *address* can be any type of addressing method used in sed, such as a single line number, a text pattern, or a range of line numbers or text patterns.

Here's an example that prints only the first two lines of a data stream to a text file:

```
$ sed '1,2w test' data6
```

The output file `test` contains only the first two lines from the input stream.

Reading from a File

You've already seen how to insert and append text into a data stream from the sed command line. The read command (`r`) allows you to insert data contained in a separate file.

The format of the read command is:

```
[address]r filename
```

The *filename* parameter specifies either an absolute or relative pathname for the file that contains the data. You can't use a range of addresses for the read command. You can specify only a single line number or text pattern address. The sed editor inserts the text from the file after the address.

```
$ sed '3r data' data2
```

The sed editor inserts the complete text from the data file into the data2 file, starting at line 3 of the data2 file.

The gawk Program

The gawk program is the GNU version of the original awk program in Unix. The awk program takes stream editing one step further than the sed editor by providing a programming language instead of just editor commands. This section describes the basics of the gawk program as a quick reference to its abilities.

The gawk Command Format

The basic format of the gawk program is as follows:

```
gawk options program file
```

[Table B.2](#) shows the options available with the gawk program.

[Table B.2](#) The gawk Options

Option	Description
<code>-F fs</code>	Specify a file separator for delineating data fields in a line.
<code>-f file</code>	Specify a file name to read the program from.
<code>-v</code> <code>var=value</code>	Define a variable and default value used in the gawk program.
<code>-mf N</code>	Specify the maximum number of fields to process in the data file.
<code>-mr N</code>	Specify the maximum record size in the data file.
<code>-W</code> <code>keyword</code>	Specify the compatibility mode or warning level for gawk. Use the <code>help</code> option to list all the available keywords.

The command line options provide an easy way to customize features in the gawk program.

Using gawk

You can use gawk either directly from the command line or from within your shell scripts. This section demonstrates how to use the gawk program and how to enter scripts for gawk to process.

Reading the Program Script from the

Command Line

A gawk program script is defined by an opening and closing brace. You must place script commands between the two braces. Because the gawk command line assumes that the script is a single text string, you must also enclose your script in single quotation marks. Here's an example of a simple gawk program script specified on the command line:

```
$ gawk '{print $1}'
```

This script will display the first data field in every line of the input stream.

Using Multiple Commands in the Program Script

A programming language wouldn't be very useful if you could execute only one command. The gawk programming language allows you to combine commands into a normal program. To use multiple commands in the program script specified on the command line, just place a semicolon between each command:

```
$ echo "My name is Rich" | gawk '{ $4="Dave"; print $0 }'  
My name is Dave  
$
```

The script performs two commands: First it replaces the fourth data field with a different value; then it displays the entire data line in the stream.

Reading the Program from a File

As with the sed editor, the gawk editor allows you to store your programs in a file and refer to them in the command line:

```
$ cat script2
{ print $5 "'s userid is " $1 }
$ gawk -F: -f script2 /etc/passwd
```

The gawk program processes all of the commands specified in the file on the input stream data.

Running Scripts Before Processing Data

The gawk program also allows you to specify when the program script is run. By default, gawk reads a line of text from the input, and then executes the program script on the data in the line of text. Sometimes you may need to run a script before processing data, such as to create a header section for a report. To do that, you use the `BEGIN` keyword. This forces gawk to execute the program script specified after the `BEGIN` keyword before reading the data:

```
$ gawk 'BEGIN {print "This is a test report"}'
This is a test report
$
```

You can place any type of gawk command in the `BEGIN` section, such as commands that assign default values to variables.

Running Scripts After Processing Data

Similar to the `BEGIN` keyword, the `END` keyword allows you to specify a program script that `gawk` executes after reading the data:

```
$ gawk 'BEGIN {print "Hello World!"} {print $0} END {print  
    "byebye"}'  
Hello World!  
This is a test  
This is a test  
This is another test.  
This is another test.  
byebye  
$
```

The `gawk` program executes the code in the `BEGIN` section first, then processes any data in the input stream, and then executes the code in the `END` section.

The `gawk` Variables

The `gawk` program is more than just an editor; it's a complete programming environment. As such, there are lots of commands and features associated with `gawk`. This section shows the main features you'll need to know for programming with `gawk`.

Built-In Variables

The `gawk` program uses built-in variables to reference specific features within the program data. This section describes the `gawk` built-in variables available for you to use in your `gawk` programs, and demonstrates how to use them.

The `gawk` program defines data as records and data

fields. A *record* is a line of data (delineated by the newline characters by default), and a *data field* is a separate data element within the line (delineated by a white space character, such as a space or tab, by default).

The gawk program uses data field variables to reference data elements within each record. [Table B.3](#) describes these variables.

[Table B.3](#) The gawk Data Field and Record Variables

Variable	Description
\$0	The entire data record
\$1	The first data field in the record
\$2	The second data field in the record
\$n	The <i>n</i> th data field in the record
FIELDWIDTHS	A space-separated list of numbers defining the exact width (in spaces) of each data field
FS	Input field separator character
RS	Input record separator character
OFS	Output field separator character
ORS	Output record separator character

In addition to the field and record separator variables, gawk provides some other built-in variables to help you know what's going on with your data and extract information from the shell environment. [Table B.4](#) shows the other built-in variables in gawk.

[Table B.4](#) More gawk Built-In Variables

Variable	Description
----------	-------------

ARGC	The number of command line parameters present.
ARGIND	The index in ARGV of the current file being processed.
ARGV	An array of command line parameters.
CONVFMT	The conversion format for numbers (see the <code>printf</code> statement). The default value is <code>%.6 g</code> .
ENVIRON	An associative array of the current shell environment variables and their values.
ERRNO	The system error if an error occurs reading or closing input files.
FILENAME	The file name of the data file used for input to the gawk program.
FNR	The current record number in the data file.
IGNORECASE	If set to a non-zero value, gawk all string functions (including regular expressions); ignore the case of characters.
NF	The total number of data fields in the data file.
NR	The number of input records processed.
OFMT	The output format for displaying numbers. The default is <code>%.6 g</code> .
RLLENGTH	The length of the substring matched in the match function.
RSTART	The start index of the substring matched in the match function.

You can use the built-in variables anywhere in the gawk program script, including the BEGIN and END sections.

Assigning Variables in Scripts

Assigning values to variables in gawk programs is similar to how you assign values to variables in a shell script—using an *assignment statement*:

```
$ gawk `
> BEGIN{
```

```
> testing="This is atest"  
> print testing  
> }'  
This is a test  
$
```

Once you assign a value to a variable, you can use that variable anywhere in your gawk script.

Assigning Variables in the Command Line

You can also use the gawk command line to assign values to variables for the gawk program. This allows you to set values outside of the normal code, changing values on-the-fly. Here's an example of using a command line variable to display a specific data field in the file:

```
$ cat script1  
BEGIN{FS=","}  
{print $n}  
$ gawk -f script1 n=2 data1  
$ gawk -f script1 n=3 data1
```

This feature is a great way to process data from your shell scripts in the gawk script.

The gawk Program Features

There are a few features of the gawk program that make it handy for manipulating data, allowing you to create gawk scripts that can parse just about any type of text file, including log files.

Regular Expressions

You can use either a Basic Regular Expression (BRE) or an Extended Regular Expression (ERE) to filter which lines in the data stream the program script applies to.

When using a regular expression, the regular expression must appear before the left brace of the program script that it controls:

```
$ gawk 'BEGIN{FS=","} /test/{print $1}' data1
This is a test
$
```

The Matching Operator

The *matching operator* allows you to restrict a regular expression to a specific data field in the records. The matching operator is the tilde circumflex (^). You specify the matching operator, along with the data field variable, and the regular expression to match:

```
$1 ~ /^data/
```

This expression filters records where the first data field starts with the text data.

Mathematical Expressions

In addition to regular expressions, you can also use mathematical expressions in the matching pattern. This feature comes in handy when matching numerical values in data fields. For example, if you want to display all of the system users who belong to the root users group (group number 0), you could use this script:

```
$ gawk -F: '$4 == 0{print $1}' /etc/passwd
```

This script displays the first data field value for all lines that contain the value 0 in the fourth data field.

Structured Commands

The `gawk` program supports the following structured commands:

The `if-then-else` statement:

```
if (condition) statement1; else statement2
```

The `while` statement:

```
while (condition)
{
    statements
}
```

The `do-while` statement:

```
do {
    statements
} while (condition)
```

The `for` statement:

```
for(variableassignment; condition; iteration process)
```

This provides a wealth of programming opportunities for the `gawk` script programmer. You can write `gawk` programs that rival the functions of just about any higher-level language program.

Index

Symbols and Numbers

: (colon)

`/etc/passwd`

`IFS`

. (dot)

BRE

regular expressions

relative filepaths

.. (dot-double), relative filepaths

; (semi-colon)

`do`

functions

`IFS`

`mysql`

`sed`

shell scripts

& (ampersand)

background mode

double-parentheses symbol

`sed`

`STDERR`

&& (ampersand-double)

double-parentheses symbol

PID

*** (asterisk)**

BRE

case

executable files

filter strings

ls

regular expressions

**** (asterisk-double), double-parentheses symbol**

\ (backslash)

for

bash shell prompt character

list

regular expressions

sed

sem

shell scripts

substitution

` (backtick)

bc

shell scripts

^ (caret symbol)

anchor character

regular expressions

sed

{ } (curly brackets), ERE

\$ (dollar sign)

bash shell

environment variables

sed

user variables

\$\$ (dollar sign-double)

bash shell

dash shell

= (equals sign)

disk space monitoring

environment variables

file permissions

print

user variables

! (exclamation mark)

double-parentheses symbol

sed

substitution

> (greater than symbol)

BSD-style parameters

output redirection

sed

>> (greater than symbol-double)

output redirection

statistics report
(hash mark)
< (less than symbol)
input redirection
sed
<< (less than symbol-double)
double-parentheses symbol
inline input redirection
- (minus sign/dash)
BRE
file permissions
format specifier
-- (minus sign/dash-double)
Mutt
shift
-? (minus sign/dash-question mark), `psql`
--- (minus sign/dash-triple), permission code
| (pipe symbol)
double-parentheses symbol
ERE
regular expressions
sort
|| (pipe symbol-double), double-parentheses symbol
+ (plus sign)
BSD-style parameters
ERE

file permissions

regular expressions

? (question mark)

ERE

filter strings

mysql

read

regular expressions

user accounts

"" (quotation marks-double)

for

IFS

positional parameters

' ' (quotation marks-single)

environment variables

sed

/ (slash)

directories

regular expressions

search

substitution

[] (square brackets)

background mode

bash shell

gawk

math

processes

regular expressions

test

[[]] (square brackets-double), if-then

] (square brackets-open), functions

~ (tilde)

double-parentheses symbol

gawk

+%y%m%d (date format)

() (parentheses)

exit status

functions

regular expressions

substitution

(()) (parentheses-double)

bash shell

if-then

zsh

\!

bash shell prompt character

mysql

\#

bash shell prompt character

\\$, bash shell prompt character

\?

mysql

psql

\@ , bash shell prompt character

\[, bash shell prompt character

\] , bash shell prompt character

\$!

bash shell

dash shell

\$#

bash shell

command line parameters

dash shell

\$*

bash shell

command line parameters

dash shell

\$- , dash shell

\$?

bash shell

dash shell

exit status

mysql

\$@

bash shell

command line parameters

dash shell

\$_ , bash shell

\$> , STDERR

\$0

bash shell

dash shell

gawk

zsh

1

exit status

signal

\$1, gawk

2

exit status

signal

\$2, gawk

3, signal

9, signal

15, signal

17, signal

18, signal

19, signal

126, exit status

127, exit status

128, exit status

128+x, exit status

130, exit status

132, exit status

255, exit status

A

A, vim

a

optstring

permissions

sem

vim

-A

fsck

Mutt

mysql

psql

tar

-a

cp

dash shell

data files

fsck

Mailx

mount

Mutt

options

psql

TCP module

\a

AAAREADME

absolute filepath

-accept_all_cookies

access

active

--active

addarray

addem

ADDR

address

--add-widget

ah

alias

aliases

aliases

Align

--all

all_export

Allow SendEvents

--almost-all

Almquist, Kenneth

[:alnum:]

[:alpha:]

Alt Sends Escape

alternative screen

Alt/Numlock Modifiers

am

amarok

American Standard Code for Information Interchange (ASCII)

anacron

anchor characters

AND

and ()

-anonymous

ANSI character set

ANSWER

Appearance

append

--append

Applets

apt-cache

apt-get

aptitude

dpkg

i u

repositories

search

sudo

wildcards

aptitude dist-upgrade

aptitude full-upgrade

ARG1 != ARG2

ARG1 \$ ARG2

ARG1 % ARG2

ARG1 + ARG2

ARG1 - ARG2

ARG1 / ARG2

ARG1 < ARG2

ARG1 = ARG2

ARG1 > ARG2

ARG1 >= ARG2

ARG1 | ARG2

ARG1<= ARG2

ARGC

ARGIND

ARGV

arrays

environment variables

gawk

variables

arraydbl

Arrow keys

ASCII. See American Standard Code for Information Interchange

ash

ash shell

ASK_COUNT

asort()

asorti()

--aspect ratio

assignment statements

-assume_charset

-assume_local_charset

-assume_unrec_charset

at

atan2()

atq

atrm

-auth

--author

autoload

Automatic Indentation

aw

B

b

block devices permissions

fdisk

optstring

vmstat

-B

-b

cp

Mailx

Mutt

mysql

sort

-b default_home

background mode

\.

--backtitle

backups

badtest

basename

BASH

bash

bash shell

for

arguments

built-in commands

CLI

command line parameters

data files

directories

disk space monitoring

environment variables

files

copying

creation
deleting
information
links
listings
filtering
renaming
viewing

filesystem navigation

manual

positional parameters

program monitoring

prompt

shell scripts

shift

starting

test

BASH_ALIASES

BASH_ARCV

BASH_ARGC

BASH_ARGV

BASH_CMDS

BASH_COMMAND

BASH_ENV

BASH_EXECUTION_STRING

bashisms

BASH_LINENO
BASHOPTS
BASHPID
.bash_profile
.bashrc
BASH_REMATCH
BASH_SOURCE
BASH_SUBSHELL
BASH_VERSINFO
BASH_VERSION
BASH_XTRACEFED
batch
bc
BEGIN
--begin x y
bel
bg
bg color
bi
-bibhost
/bin
bindkey
/bin/sh
[[:blank:]]
blink
BLOB

blocking

block devices

files

permissions

block mode graphics

block preallocation

`--block-size=size`

`bold`

`-book`

Bookmarks menu bar

`Boolean`

Boolean operators

`/boot`

boot process

`bounce`

Bourne again shell

Bourne shell

Bracket Matching

`branch`

BRE. See POSIX Basic Regular Expression engine

`break`

bash shell

`while`

Break key

broken dependencies

`bsd_echo`

BSD-style parameters

buff

buffers

emacs

graphics

scrollback

--buffer-size=*SIZE*

builtin

built-in commands

bash shell

dash shell

zsh

built-in functions, [gawk](#)

built-in variables, [gawk](#)

bunzip2

-buried_news

bye

bzcat

bzip2

bzip2recover

C

c, [sendmail](#)

c

aptitude

character device permissions

fdisk

format specifier

optstring

-C

fsck

lvcreate

mysql

sort

-c

dash shell

exit

lvcreate

Mailx

Mutt

options

psql

sort

tar

TCP module

C-@

\C

\c

-c bytes

-c *cmdlist*

-c *comment*

-c *dir*

C language

-c string

C-a

cache

-cache

calendar

--calendar

caller

calling functions

--cancel-label

Capitalize

case

command line parameters

getopts

menus

user accounts

zsh

-case

case sensitivity

cat

bash shell

inputbox

read

zsh

C-b

c_bases

cd

bash shell

dash shell

zsh

CDPATH

C-e

-center

C-f

-cfg

chage

change

Change Case

Change Profile

char

Character Encoding

characters

anchor

ANSI

classes

devices

files

permissions

escape

ranges

special

charset

chdir
--check
CHECK_DIRECTORIES
checklist
--checklist
check=none
--check=quiet
chfn
chgrp
-child
chmod

bash shell
permissions

chown
chpasswd
chsh
--chunksize
C-k
--classify

Clean Indentation

cleanup
clear
--clear

Clear Scrollback & Reset

CLI. See command line interface
client/server paradigm

cm

CMD

-cmd

-cmd_log

-cmd_script

C-n

color

Color Scheme

--colors

--color=when

--cols *n*

cols#80

--column

COLUMNS

--columns *n*

COM port

COMMAND

command

--command

commands. *See also specific commands*

aliases

bash shell

C language

groups, sed

values

command line interface (CLI)

bash shell

functions

gawk

MySQL

PostgreSQL

variables

command line parameters

bash shell

case

chage

dash shell

GNOME

Konsole

ls

lynx

options

tail

test

useradd

while

xterm

zsh

commands

COMP_CWORD

compgen

COMP_KEY

compl ()
complete
COMP_LINE
COMP_POINT
COMPREPLY
compress
COMP_TYPE
COMP_WORDBREAKS
COMP_WORDS
COM_WORDBREAKS
--concatenate

conditions

CONFIG_FILE
configure

Configure Editor

configure sysstat
connect
_connect_timeout
CONT

--context
continue
--continuous

control codes

terminal emulation
terminfo database

CONVFMT

`_cookie_file`

`_cookies`

`_cookie_save`

copy

copy and paste

`emacs`

`vim`

Copy as HTML

`-core`

core distributions

`corutils`

`cos()`

counting parameters

`cp`

`C-p`

`%CPU`

`cr`

`C-r`

`-crawl`

`--create`

`--create-rc`

`cron`

`Delete_User`

statistics report

`crontab`

`--cr-wrap`

cs

C-s

csch_junkie_history

csch_junkie_loops

csch_junkie_quotes

csch_nullcmd

C-Spacebar

csr

cub

cubl

cud

cudl

cuf

cuf1

--cumulative

cup

cURL

curl

Current Line

-curses_pads

Cut

cuu

cuul

C-v

C-x 2

C-x 3

C-x b
C-x C-b
C-x C-c
C-x C-f
C-x C-s
C-x o
C-y
C-z

D

D, sendmail

d

directory permissions

fdisk

format specifier

optstring

sed

sem

d\$

-D

Mutt

mysql

useradd

-d

chage

cp

if

mktemp

options

psql

sort

tar

\d

-d file

-d *home_dir*

da*

Daily_Archive

dash shell

built-in commands

command line parameters

echo

environment variables

function

positional parameters

shell scripts

test

user-defined environment variables

user-defined variables

data

archives

compression

constraints

emacs

fields

files

-a

bash shell

formatting

input

file descriptors

shell scripts

output

file descriptors

shell scripts

variables

presentation

search

sorting

databases

MySQL

PostgreSQL

problem-tracking

shell scripts

tables

DATE

Date

date

bash shell

if

dbl

dbl \$value

--dbname

dd

deb

Debian

local startup file locations

MySQL

PMS

run level

deb-src

-debug_partial

DEC. See Digital Equipment Corporation

declare

--default-item

--defaultno

deferred

-delay

DELETE

delete

bash shell

directories

read

sed

sem

--delete

Delete key

Delete_User

delimiter

dependencies

broken

--dereference

DESCRIBE

Deselect

--deselect

desktop environment

DESTINATION

destination

/dev

DEVICE

devices

names

permissions

/dev/null

/dev/sdcl

/dev/sdx

df

dialog

dialog package

Dickey, Thomas E.

--dictionary-order

--diff

digiKam

[[[:digit:]]

Digital Equipment Corporation (DEC)

directories

for

bash shell

cron

delete

files, counting

ls

permissions

temporary

virtual

wildcards

--directory

dirs

DIRSTACK

disable

--disable-factory

disk space monitoring

bash shell

diskmail

DISKS_TO_MONITOR

disown

-display

-display_charset

distributions

core

LiveCD

specialized

distrubution_name

do

Document Statistics

Documents command

dolphin

Domain Name System (DNS)

domaintable

done

-dont_wrap_pre

Doublesized characters

do-while

dpkg

drag-and-drop

\dt

du

\du

dumb terminals

-dump

dw

dyne:bolic

E

e

format specifier

-E

chage

dash shell

mysql

psql

-e

dash shell

file comparison

Mailx

Mutt

mysql

options

psql

\e

-e command

-e *expiration_date*

-e *expire_date*

-e file

-e *script*

echo

bash shell

dash shell

disk space monitoring

file descriptors

Mailx

menus

-n

STDERR

while

zsh

--echo-all

--echo-hidden

--echo-queries

ed

edit

Edit command

emacs

gedit

KWrite

Edit menu bar

GNOME terminal emulation

Konsole

_editor

editors

emacs

GNOME

KDE

vim

-ef

ego

egrep

8-bit controls

el

elif

ell

EMACS

emacs

- buffers

- console mode

- copy and paste

- data

- graphical windows

- menu bar

- search and replace

- X Windows

-emacskeys

e-mail

- parsing, regular expressions

- STDERR

- STDOUT

empid

emulate

-en

enable

-enable_scrollback

enacs

Encoding

--encoding

END

end file

end-of-file (EOF)

--entry

ENVIRON

environment variables

- arrays

- bash shell

- Bourne shell

- command aliases

- dash shell

- defaults

- global

- local

- locating

- Lynx

- removing

- setting

- shell scripts

EOF. See end-of-file

epiphany

equal sign

ERE. See POSIX Extended Regular Expression engine

err_exit

ERRNO

error

--error

-error_files

err_return

--escape

escape characters

/etc

/etc/apt/sources.list

/etc/bashrc

/etc/fstab

/etc/group

/etc/init.d

/etc/init.d/boot.local

/etc/init.d/rc.local

/etc/inittab

/etc/passwd

: (colon)

chfn

IFS

/etc/profile

etc/rc#

/etc/rc.d/rc.local

/etc/rc.local

/etc/rcX.d

/etc/shadow

/etc/skel

/etc/urpmi/skip.list

EUID

eval

eval_lineno

evince

Evolution

exec

date

file descriptors

STDINPUT

STDOUT

while

zsh

-exec

executable files

exit

bash shell

-c

dash shell

mysql

zsh

exit status

dialog

dialog widgets

--exit-label

exp()

--expanded

export

expr

(EXPRESSION)

ext

ext2

ext3

ext4

extend command

extended filesystem

Extensible Markup Language (XML)

extents

--extents

External Tools

--extract

F

F

f

-F

mount

Mutt

psql

-f

cp

dash shell

file comparison

funclib

gawk

if

mount

Mutt

mysql

options

psql

sort

TCP module

-f *file*

sed

-F fs

-f *inactive*

-f inactive_days

false

FAT16

FAT32

fb font

fc

FCEDIT

FD

fdisk

Fedora

ext4

local startup file locations

yum

fg

fg color

fgrep

field separators

IFS

OFS

--field-separator

FIELDWIDTHS

FIGNORE

file

bash shell

Mutt

\$file

--file

files. *See also specific file types*

block devices

character devices

comparisons, `test`

copying, bash shell

creation, bash shell

deleting, bash shell

descriptors

choosing

-d

data input

data output

echo

exec

read

read/write

redirecting

snapshot reports

STDIN

STDOUT

gawk

globbing

information, bash shell

links, bash shell

listings

filtering

loops

network devices

permissions

decoding

defaults

groups

ownership

security

symbols

renaming, bash shell

sed

sharing

viewing, bash shell

File Browser Pane

File command

emacs

gedit

KWrite

File menu bar

GNOME terminal emulation

Konsole

File Saving

`file1 -nt file2`

`file1 -ot file2`

`FILE_LIST`

`FILENAME`

`FILE_NO`

filepath

`--file-selection`

filesystems

creating

journaling

management

kernel

LVM

navigation, bash shell

partitions

--file-type

filtering

Find

find

Find Next

Find Previous

Find Selected

Find Selected Backwards

Find_Problem

finger

FIXED_DATE

fixed_date

float

floating-point solutions

flush

fluxox

fn font

FNR

folding marker

fonts

for

bash shell

break

C language
directories

do

done

gawk

IFS

multiple variables

sort

variables

wildcards

zsh

-force_empty_hrefless_a

-force_html

-force_secure

--forest

form

--format *format*

format specifiers

format strings

-forms_options

fpath

fragmentation

free

-from

FS

fsck

fselect

-ftp

func1

funclib

FUNCNAME

function

functions

- array variables

- CLI

- libraries

- menus

- parameters

- recursion

- shell scripts

- values

- variables

- zsh

Function keys

function_argzero

functions

fwwm

fwwm95

fw font

G

G

sed

vim

g

format specifier

group permissions

sed

substitution

-G

file comparison

groups

mysql

-g

groups

sort

\G

\g

-G file

-G *group*

-g *group*

-G *grplist*

-g *grplist*

-g *initial_group*

gauge

gawk

for

`[]` (square brackets)

advanced commands

arrays

`BEGIN`

built-in functions

built-in variables

CLI

disk space monitoring

`do-while`

`END`

`-f`

files

`if`

matching operator

math

multiple commands

options

`print`

`printf`

quick reference

regular expressions

`--re-interval`

string functions

structured commands

time functions

user-defined functions

user-defined variables
variables

`while`

zombie processes

gcalc-tool

`gdialog`

gedit

`--general-numeric-sort`

`gensub()`

Gentoo

`--geometry=`

`get_answer`

`get_data`

`getln`

`getopt`

`getopts`

bash shell

dash shell

zsh

`gg`

global environment variables

global variables

`global_export`

`global_rcs`

`GLOBIGNORE`

GNOME. See GNU Network Object Model Environment

gnome-nettool

gnome-panel

gnome-terminal

GNU long parameters

GNU Network Object Model Environment (GNOME)

- applications

- command line parameters

- editors

- e-mail

- graphical windows

- terminal emulation

 - menu bar

 - tabs

- Ubuntu

GNU utilities

go

Go to Line

GRANT

grant

graphics

- block mode

- buffering

- color

- terminal emulation

- vector

graphical desktops

graphical virtual consoles

graphical windows

GNOME

KDE

menus

grep

groups

permissions

regular expressions

role

sed

SGID

VG

--Group *grplist*

--group *grplist*

groupadd

--group-directories-first

groupmod

GROUPS

grub2-theme-mint

gsub()

gunzip

gzcat

gzip

H

H

sed

sendmail

h

sed

vim

-H

Mutt

mysql

-h

Mutt

mysql

options

psql

\H

\h

hard drives. *See also* disk space monitoring

hard links

hardware management

hash

hc color

head

-head

--headers

height

help

-help

--help

Help command

Help menu bar

--help-button

--help-label

--help-status

hidden files

-hiddenlinks

Hide VT Window

--hide-menubar

histchars

HISTCMD

HISTCONTROL

HISTFILE

HISTFILESIZE

HISTIGNORE

-historical

history

HISTSIZE

hold space

HOME

useradd

userdel

home

\$HOME

disk space monitoring

temporary files

/home

\$HOME/bin

\$HOME/mail

\$HOME/.my.cnf

-homepage

\$HOME/.zshrc

--host

HOSTFILE

HOSTNAME

HOSTTYPE

Hourly_Archie

hpfs

ht

HTML

sed

statistics report

hts

HTTP headers

--human-readable

HUP

I

I

i

format specifier

sem

-I

chage

dash shell

lvcreate

mount

-i

cp

dash shell

lvcreate

Mailx

mysql

options

sort

zsh

i u

Id

id

identified by

ID_NUMBER

id_number

if

-d

date

-f

gawk

IFS. **See** internal field separators

if-then

(()) (parentheses-double)

[][] (square brackets-double)

break

continue

read

test

if-then-else

--ignore

--ignore-backups

IGNORECASE

--ignore-case

IGNOREEOF

--ignore-leading-blanks

-image_links

incoming

incremental search

ind

Indent

index()

-index

index STRING CHARS

index values

--info

infobox

infocmp

init process

inline input redirection

inner loops

--inode

inode numbers

input redirection

inputbox

--inputbox

--input-fd

inputmenu

INPUTRC

--insecure

INSERT

insert

Insert Date/Time

insert mode

install

install.cf

INT

int ()

int

integer

interactive shell

internal field separators (IFS)

International Organization for Standardization (ISO)

Internet

cURL

Lynx

zsh

Internet Message Access Protocol (IMAP)

intervals

ipcs

-ismap

ISO. See International Organization for Standardization

.iso

iso9660

--item-help

J

J, vim

j

vim

-j

\j

JFS. See Journaled File System

jfs

jfs_mkfs

jobs

job numbers

Joe's Window Manager (JWM)

Join Lines

Journaled File System (JFS)

journaling

JWM. See Joe's Window Manager

K

κ, sendmail

k

vim

-k

sort

useradd

K Desktop Environment (KDE)

applications

editors

e-mail

graphical Windows

K menu

κ *sort*

K3b

ka3

Kaffeine

ka1

Kate

kb2

kbs

kc1

kc3

kcubl

kcud1

kcuf1

kcuul

KDE. See K Desktop Environment

kdiallog

--keep-window

kent

kernel

- device drivers

- filesystem management

- hardware management

- memory management

- modules

- software program management

- virtual memory

--key

keyboard

- terminal emulation

- xterm

--keytab file

--keytabs

kf0

kf1

kf10

KILL

kill

killall

KMail

Knoppix

Koffice

Konqueror

Konsole

- command line parameters

- menu bar

- profiles

- tabbed windows

Kontakt

Kopete

korn

ksh_arrays

ksh_autoload

ksh_option_print

ksh_typeset

KWrite

kwrite

L

L

l

fdisk

links permissions

print

vim

-L

lvcreate

TCP module

usermod

-l

dash shell

jobs

lvcreate

Mailx

mount

options

psql

TCP module

usermod

zsh

\l

-L *label*

Lam, Savio

LAMP. See Linux-Apache-MySQL-PHP

LANG

Latin-1

LC_ALL

LC_COLLATE

LC_CTYPE

LC_MESSAGES

LC_NUMERIC

length()

length STRING

less

let

bash shell

zsh

lf filename

/lib

libraries

LIMIT

limit

--line

line addressing

Line Drawing Characters

Line Numbers

LINENO

LINES

--lines *n*

lines#24

link

-link

links

Linux console

Linux-Apache-MySQL-PHP (LAMP)

linuxmint.png

list

\$list

--list

psql

tar

lists

LiveCD distributions

ln

local

bash shell

zsh

local environment variables

local installation

local variables

-localhost

local_options

local_traps

-locexec

log()

log

Log to file

logfiles

logging messages

Logical Volume Manager (LVM)

- filesystem management

- mirroring

- modifying

- snapshots

- striping

logical volumes

login

Login Roles

login shell

logout

loop

loops. *See also specific loop operators*

- controlling

- files

- nested

- output processing

`[[:lower:]]`

Lowercase

ls

- * (asterisk)

- bash shell

- command line parameters

- directories

- filtering

- regular expressions

STDOUT

--ls

LS_COLORS

lshift()

lsof

lvcreate

lvdisplay

lvextend

LVM. See Logical Volume Manager

LVM1

LVM2

lvreduce

lvtest

Lynx

environment variables

output capturing

lynx

LYNX_CFG

LYNX_LSS

LYNX_SAVE_SPACE

M

M, sendmail

m

fdisk

regular expression interval

-M

sort

useradd

-m

chage

dash shell

lvcreate

Mutt

sort

useradd

M-<

M->

M-a

MACHTYPE

MAIL

mail

Mail Delivery Agent (MDA)

Mail Transfer Agent (MTA)

Mail User Agent (MUA)

MAILCHECK

Maildir-style mailboxes

maildrop

mailertable

MAILPATH

mailq

MAIL_TO

Mailx

main.cf

make

man

Mandriva

master.cf

match()

match STRING REGEXP

matching operator

math

gawk

shell scripts

zsh

Mauelshagen, Heinz

maverick

Maverick Mercat

--max-input

mb

M-b

M-backspace

M-d

M-e

media

/media

%MEM

MemFree

meminfo

memory

management

pages

MemTotal

menu

while

--menu

menus

functions

graphical windows

select

shell scripts

menu bar

emacs

GNOME terminal emulation

Konsole

--merge

messages

area

IMAP

logging

shell scripts

Meta sends Escape

M-f

-mf

-mf N
-mime_header
-minimal
minix
--minor

Mint

mirroring

M-k
mkdir
mke2fs
mkefs
mkfs.ext3
mkfs.ext4
mkfs.xfs
mkreiserfs
mktemp
mktemp -d
mktime()
m.n
/mnt
mode

mode table
Modelines
modules
 kernel
 zsh

- modulo
- monitoring
 - disk space
 - bash shell
- processes
- program
- real-time process
- X Windows

--month-sort

more

mount

mount

- bash shell
- hard drives

mount points

mounting media

moving

mqueue

-mr N

ms color

msdos

msgbox

--msgbox

msgr

multios

multiple commands

gawk

sed

shell scripts

while

multiple variables, for

Mutt

mutt

mv

M-v

\$myarray

myfuncs

myprint()

mypro?

myprog

MySQL

CLI

mysql

\$?

-e

mysql INSERT

mysql UPDATE

N

N

sed

n

fdisk

network device permissions

sed

\$n

dash shell

gawk

-N

fsck

mysql

-n

disk space monitoring

echo

jobs

lvcreate

Mailx

mount

Mutt

options

print

read

sed

sort

substitution

test string comparison

useradd

\n

-n lines

N *namelist*

-n *namelist*

-n str1

n1 eq n2

n1 ge n2

n1 -le n2

n1 -ne n2

n1-gt n2

n1-lt n2

NAME

name

--name

name name

--name name

nautilus

ncp

nested loops

-nested_tables

network

devices

files

permissions

zsh

Network File System

[New File]

newaliases
-newschunksize
-newsmaxchunk
next
NF
nfs
nftfs
NI
nice
\nnn
NNTPSERVER
--no-align
-nobold
-nobrowse
--nocancel
-nocc
--noclose
--no-collapse
-nocolor
NODE

nodes

-noexec
-nofilereferer
--noframe
--no-group
--no-headers

--nohost

nohup

nohup.out

--no-kill

--no-label

-nolist

-nolog

--nomenubar

non-incremental search

non-interactive shell

nopager

-nopause

-noprint

--nopsqlrc

-noredir

-noreferer

--noresize

-norestarting_sigwinch

-noreverse

normal mode, vim

--no-shadow

-nostatus

not null

--notabbar

note

--notification

-nounderline

nowarning

--noxft

NR

-nt

null files

NULLCMD

num _G

-number_fields

-number_links

numeric comparisons, [test](#)

numeric line addressing, [sed](#)

--numeric-sort

--numeric-uid-gid

-nw

O

O

[sendmail](#)

o

[fdisk](#)

[file permissions](#)

[format specifier](#)

-O

[file comparison](#)

[mount](#)

-o

mount

mysql

options

psql

sort

-O file

o *format*

o *format*

-o *format*

-o *format*

o *order*

octal mode

octal_zeroes

OFMI

OFMT

OFS

--ok-label

Old Function keys

OLDPWD

Open File Manager

open source software (OSS)

openSUSE

local startup file locations

MySQL

zypper

/opt

OPTARG

OPTERR

OPTIND

options

chmod

combining

command line parameters

gawk

standardizing

tar

values

zsh

Options command

optstring

OR

or()

o+r

ORDER BY

Ordered mode

ORS

OSS. See open source software

OSTYPE

outchars

outer loops

output

capturing

Lynx

statistics report

xterm

file descriptors

shell scripts

redirection

variables

--output

psql

output field separator. *See* `OFS`

output record separator. *See* `ORS`

--output-fd

Overwrite Mode

P

`P`, `sendmail`

`p`

`fdisk`

`print`

`sed`

`sem`

`substitution`

`-p`

`cp`

`jobs`

lvcreate

Mutt

mysql

PID

psql

read

tar

usermod

zsh

\P

\p

-p *num*

-p *passwd*

p *pidlist*

-p *pidlist*

Package Management System (PMS)

Debian

package_type_list

Packed Font

packets

PageDown

pager

PageUp

parameters. See also specific parameter types

functions

-partial

-partial_thres

partitions

filesystems

PV

passwd

bash shell

--password

passwords

/etc/shadow

passwordbox

passwordform

Paste

PATH

setting

/usr/sbin

pattern space

pause

-pauth

PCLinuxOS

permissions

chmod

devices

files

groups

users

--persistent

phone number validation

physical volumes (PV)

pickup

PID. See process ID

--pid *pidlist*

-pid=PID

pipe

pipes

PIPESTATUS

plain text

plug-ins

PMS. See Package Management System

POP. See Post Office Protocol

popd

-popup

positional parameters

POSIX Basic Regular Expression engine (BRE)

. (dot)

* (asterisk)

anchor characters

character classes

character ranges

plain text

special characters

POSIX Extended Regular Expression engine (ERE)

{ } (curly brackets)

| (pipe symbol)

+ (plus sign)

? (question mark)

posix_builtins

POSIXLY_CORRECT

--post

Post Office Protocol (POP)

-post_data

postdrop

Postfix

PostgreSQL

CLI

PPID

--ppid *pidlist*

PR

prec

Preferences dialog

-preparsed

-prettysrc

PRI

primary key

primary partition

print

gawk

sed

zsh

-print

[:print:]

Print window

printenv

printf

--print-maxsize

--print-size

problem-logger

problem-tracking databases

Problem_Trek

prob_solutions

prob_symptoms

proc

process ID (PID)

&&

while

process ID (PID)

process_answer

processes

monitoring

real-time

signals

status

stopping

user accounts

virtual memory

zombie

/proc/meminfo

--profile file

profiles

--profiles

program monitoring

Program shortcuts

--progress

prompt

PROMPT_COMMAND

PROMPT_DIRTRIM

PROTOCOL_PROXY

ps

PS1

PS2

PS3

PS4

--pset

-pseudo_inlines

psql

public schema

[:punct:]

Puppy Linux

pushd

pushln

PV. See physical volumes

pvcreate

pvdiskplay

PWD

pwd

bash shell

dash shell

zsh

Python Console

Q

q

fdisk

vim

q!

-Q

-q

mysql

options

psql

\q

mysql

psql

-qa

qmgr

--question

Quick Open

--quiet

QUIT

quit

R

R

r

sed

sem

vmstat

-R

chown

cp

fsck

lvcreate

Mutt

psql

sort

-r

cp

file comparison

fsck

gawk

jobs

lvcreate

mount

mysql

options

tar

useradd

userdel

zsh

r--

\R

\r

r *char*

-r file

R *text*

radiolist

--radiolist

rand()

RANDOM

--random-sort

--random-source=*FILE*

-raw

rc

rcs

read

? (question mark)

bash shell

cat

dash shell

delete

file descriptors

if-then

menus

-n

-p

pipe

-s

sed

STDIN

-t

user accounts

values

while

zsh

Read Only Mode

--readahead

README

READNULLCMD

readonly

bash shell

dash shell

zsh

read/write file descriptors

-realm

real-time process monitoring

recipients

Record_Problem

records

--record-separator

RECORDS_EXIST

recursion

--recursive

Red Hat

run level

yum

Red Hat Package Management (RPM)

MySQL

Redirect to printer

Redo

Redraw Window

REG

REGEXP

--regionsize

regular expressions

*** (asterisk)**

case sensitivity

counting directory files

e-mail parsing

engine

gawk

groups

phone number validation

sed

wildcards

rehash

--re-interval

Reiser, Hans

ReiserFS

relative filepath

relay-domains

-reload

removable media

renice

Repeat key

Replace

replace-string

REPLY

report_date

REPORT_FILE

repositories

aptitude

urpm

versions

yum

zypper

RES

-restrictions

-resubmit_posts

\$result

return

Return key

rev

reverse

--reverse

ri

rich

Right margin

RLENGTH

-rlogin

rm

rmacs

rmam

rmdir

rmkx

rmso

rmul

ro

--role+

rolling windows

root

file owners

MySQL

problem-solver

test

/root

root drive

--rows *n*

RPM. See Red Hat Package Management

rpm

rpmfusion.org

rpm.list

RS

rs2

rshift()

RSS

RSTART

run level

Debian

init process

Red Hat

rv

rw

rwX

r-x

S

S

sendmail

s

fdisk

format specifier

-S

cp

mount

mysql

psql

sort

-s

cp

dash shell

disk space monitoring

file comparison

fsck

jobs

lvcreate

Mailx

Mutt

mysql

options

psql

read

soft links

sort

STDOUT

zsh

\s

-s file

-s sec

-s *sesslist*

-s *shell*

safe-upgrade

sb

/sbin

sc

--scale

schema

--schema name

scratch area

screen scraping

script. See *also* shell scripts

startup

Scroll Lock key

scroll region

scrollback buffer

Scrollback menu bar

sda

search, data

search

search and replace

search and substitute

Search command

SECONDS

Secure Keyboard
security
changing settings

sed

& (ampersand)
! (exclamation mark)
advanced commands

append

branch

change

changing flow
commands

groups
multiple
negating

d

delete

disk space monitoring
double spacing lines

equal sign

files
hold space
HTML

insert

line addressing

list

-n

next

numbering lines

numeric line addressing

p

print

printing last lines

quick guide

r

read

regular expressions

reverse

shell scripts

STDIN

STDOUT

substitution

t

test

text pattern filters

transform

utilities

wrappers

SEGV

SELECT

select

menus

zsh

Select All

-selective

self-containment

sem

a

\ (backslash)

d

delete

i

p

r

w

y

sendmail

sendmail.cf

sendmail.ct

sendmail.cw

--separate-output

--separate-widget

--separator

set

bash shell

dash shell

getopt

zsh

set group id (SGID)

set user id (SUID)

setopt

Settings command

Settings menu bar

SGL. *See* Silicon Graphics Incorporated

SGID. *See* set group id

sgr

sgr0

--shadow

shadow files

sharing

files

memory pages

shell. *See also specific shells*

terminal emulation

SHELL

shell scripts

advanced

background mode

bash shell

controls

dash shell

data

input

output
databases
dialog
environment variables
exiting
functions
getopt
graphical desktops
input redirection
jobs
math
menus
messages
multiple commands
output redirection
pipes
problem-tracking databases
sed
snapshot report
statistics report
STDERR
suppressing command output
utilities
variables
zenith
zsh

SHELLOPTS

sh_file_expansion

shift

-- (dash-double)

bash shell

dash shell

zsh

SHLVL

sh_nullcmd

shopt

sh_option_letters

-short_url

SHOW

Show Alternate Screen

Show Menu Bar

SHOW TABLES

Show Tek Window

-show_cursor

-show_dquotes

--show-menubar

showq

-show_rate

SHR

sh_word_split

si

--si

--sid *sislist*

SIGCONT

SIGHUP

SIGINT

SIGKILL

signals

SIGQUIT

SIGTERM

SIGTSTP

Silicon Graphics Incorporated (SGI)

Simple Mail Transfer Protocol (SMTP)

SimplyMEPIS

sin()

--single-line

--single-quoted

--single-step

SIZE

--size

Slackware

Slax

--sleep

smacs

smam

smb

smkx

smso

smtp

smtpd

smul

--snapshot

snapshots

LVM

report

Snippets

SNTP. See Simple Mail Transfer Protocol

so

soft links

software

installation

program management

--sorry

Sort

sort

for

bash shell

rpm

test

--sort *order*

--sort=extension

--sort=none

--sort=size

--sort=time

--sort=version

sound juicer

source

mysql

zsh

-source

source code

sources.list

[:space:]

special characters, BRE

Spell Checker

Spellcheck Selection

Spelling

Spelling (from cursor)

split()

sprintf()

sqrt()

sshd

SSL_CERT_DIR

SSL_CERT_FILE

--stable

-stack-dump

-startfile_ok

startup script

--startup-id=

STAT

stat

statistics report

status

STDERR

& (ampersand)

/dev/null

dialog

dialog widgets

echo

e-mail

nohup

null files

shell scripts

testerror

--stderr

STDIN

file descriptors

gawk

read

sed

testfile

-stdin

--stdin

STDINPUT

STDOUT

cURL

databases

e-mail

`exec`

file descriptors

`gawk`

`ls`

`-n`

`nohup`

`-s`

`sed`

snapshot reports

`STDERR`

`substitution`

temporary files

`--stdout`

sticky bit

STIME

`STOP`

`str1 != str2`

`str1 < str2`

`str1 = str2`

`str1 > str2`

stream editors

`strftime()`

string

comparisons, test

functions, gawk

String : REGEXP

--stripes

striping

--stripsize

structured commands

gawk

zsh

stub functions

su-

sub()

substitution

escape characters

N

sed

substr()

substr STRING POS LENGTH

sudo

SUID. See set user id

Sun Function keys

suspend

SW

swap space

swapping out

SwapTotal

Switch to Tek Window

swpd

sy

Synaptic Package Manager

sysstat

system

system accounts

System V init process

sys_time()

sysv

SZ

T

T

t

fdisk

sed

-T

fsck

mysql

psql

sort

-t

filesystems

fsck

mysql

psql

read

sort

tar

TCP module

zsh

\T

\t

t. --title=

-T title

t *ttyst*

-t *ttyst*

--tab

Tab Stops

tabbed windows

--tab-correct

<table>

tables

cron

databases

--table-attr

--tab-len

tabs, GNOME terminal emulation

Tabs menu bar

--tab-with-profile=

tac

Tag List

tags

-tagsoup

tail

tailbox

tailboxbg

tar

tasks

tb

tbc

TCP module

tcsh

tee

Tektronix

-telnet

\$temp

\$tempfile

temporary directories

temporary files

\$HOME

--temporary-direction=*DIR*

\$tempt

TERM

-term

Termcap Function keys

terminal emulation

control codes

GNOME

graphics

keyboard

Konsole

Linux Console

shell

terminfo database

Terminal menu bar

terminfo database

control codes

test

bash shell

command line parameters

dash shell

file comparisons

if-then

numeric comparisons

root

sed

sort

string comparisons

zsh

\$test

test command

testerror

testfile

Text

text

text pattern filters

text virtual consoles

Text Wrapping

textbox

--textbox

--text-info

TIME

Time

TIME+

time functions

timebox

TIMEFORMAT

--timeout

times

Timestamp

--title

-tlog

TMOUT

/tmp

TMPDIR

-tna

+ TOKEN

tolower()

tomboy

Tools command

top

Torvalds, Linus

totem

touch

touch top

toupper()

-trace

-trace_mask

transform

trap

traps_async

-traversal

--trim

-trim_input_fields

trivialrewrite

true

TrueType Fonts

TSTP

TTY

--tty *ttylist*

ttyctl

ttys

--tuples-only

TYPE

type

--type type
--types
typeset
typeset_silent

U

u
fdisk
user permissions

vim
-U
mysql
psql
usermod

-u
cp
dash shell

find
mysql
sort
tar
zsh

\u
-u *uid*
U *userlist*
-U *userlist*

-u *userlist*

-U *uuid*

Ubuntu

/etc/init.d

ext4

GNOME

local startup file locations

Maverick Mercat

MySQL

UCE. See unsolicited commercial e-mail

ufs

UID. See user ID

\$UID

ulimit

umask

umount

umsdos

unalias

-underline_links

-underscore

Undo

unfunction

unhash

Unicode

Uniform Resource Identifier (URI)

Unindent

--unique

Unix-style parameters

unlimit

unmount

unset

bash shell

dash shell

zsh

unsetopt

unsolicited commercial e-mail (UCE)

until

break

continue

while

unzip

Update_Problem

[:upper:]

Uppercase

Upstart init process

uptime

URI. See Uniform Resource Identifier

urpm

Mandriva

repositories

urpme package_name

urpmi --auto-update --update

urpmi --clean
urpmi package_name
urpmi.addmedia path_name
urpmq --list-media
us
USE
use
--use-factory
-use_mouse
USER
user
\$USER

users

accounts
 processes
adding
modifying
permissions
removing
user ID (UID)
 bash shell
--User *userlist*
--user *userlist*
user variables
USER_ACCOUNT_RECORD
useradd

-useragent

user-defined environment variables

user-defined functions

user-defined variables

userdel

userid

usermod

--username

/usr

/usr/sbin

UTF-8 Fonts

UTF-8 Titles

utilities

GNU

sed

shell scripts

u-x

V

V

v

fdisk

-V

dash shell

fsck

psql

-v

cp

dash shell

gawk

Mailx

mount

Mutt

mysql

options

psql

tar

TCP module

\V

\v

-v var=value

val++

++val

val--

--val

-validate

\$value

values

commands

functions

options

read

/var

varchar

--variable

variables. See *a/so* environment variables

for

arrays

CLI

dash shell

data fields

data output

functions

gawk

global

lists

local

shell scripts

user

var/mail

/var/spool/mail

vector graphics

verbose

-verbose

-version

--version

versions, repositories

vfat

VFS. See Virtual File System

VG. See volume groups

`vgchange`

`vgcreate`

`vgdisplay`

`vgextend`

`vgreduce`

`vgremove`

video card

View command

`gedit`

`KWrite`

View menu bar

GNOME terminal emulation

Konsole

`-vikeys`

`vim`

copy and paste

search and substitute

VIRT

virtual consoles

virtual directory

Virtual File System (VFS)

virtual memory

`virtusertable`

`--visit`

visual mode

vmstat

volume groups (VG)

VSZ

VT commands

VT features

VT fonts menu

VT options menu

--vt sz CxL

vt#3

VT220 keyboard

W

W

w

fdisk

sed

sem

substitution

-W

chage

gawk

mount

psql

-w

file comparison

mysql

-w-, permission code

\W

mysql

\w

mysql

-w file

w *filename*

-W keyword

wa

wait

bash shell

dash shell

zsh

--warning

--warningcontinuecancel

warnings

--warningyesno

--warningyesnocancel

wc

WCHAN

whence

WHERE

where

which

bash shell

Mailx

zsh

while

break

command line parameters

continue

echo

exec

gawk

getopts

menus

multiple commands

PID

read

until

user accounts

zsh

who

widget

widgets

dialog package

kdialog

zenity

width

-width

--width *n*

wildcards

aptitude

directories

regular expressions

rm

--window

windows. *See also* graphical windows

rolling

tabbed

--window-with-profile=

-with_backspaces

Word Wrap Document

--workdir dir

--working-directory=

wq

wrappers

Writeback mode

WWW_HOME

-wx

X

X

format specifier

x

fdisk

format specifier

sed

vim

-X, psql

-x

cp

dash shell

file comparison

Mutt

mysql

options

psql

tar

--x, permission code

x < y

x <= y

x == y

x > y

X event commands

-x file

X Windows

emacs

virtual consoles

x >=y

Xandros

xen1

xfce

XFree86

XFS

XFS filesystem

XML. See Extensible Markup Language

xon

xor()

X.org

xterm

- command line parameters

- keyboard settings

- main menu

- output capturing

- VT options menu

- X event commands

xtrace

Y

y

- sem

-y

- fsck

Mutt

options

--yes-label

yesno

--yesno

--yesnocancel

+%y%m%d

yum

Fedora
Red Hat
repositories

yum localinstall

Z

Z

z

-Z

lvcreate

Mutt

-z

ANSWER

Mutt

sort

tar

test string comparison

\z

z shell (zsh)

- built-in commands
- command line parameters
- functions
- mathematical operations

modules

options

shell scripts

structured commands

-z str1

zcompile

zenity

--zero

--zero-terminated

zftp

zip

zipcloak

zipnote

zipsplit

zmodload

\$HOME/.zshrc

zsh

zombie processes

--zoom=

zsh. See z shell

zsh

zsh/cap

zsh/clone

zsh/compctl

zsh/complete

zsh/complist

zsh/computil
zsh/datetime
zsh/deltochar
zsh/files
zsh/mapfile
zsh/mathfunc
zsh/net/socket
zsh/net/tcp
zsh/parameter
zsh/pcre
zsh/sched
zsh/stat
zsh/system
zsh/termcap
zsh/terminfo
zsh/zftp
zsh/zle
zsh/zleparameter
zsh/zprof
zsh/zpty
zsh/zselect
zsh/zutil
ztcp
zypper

openSUSE
repositories

zypper addrepo path_name

zypper install package_name

zypper remove_package_name

zypper repos

zypper update

zypper verify